



OXFORD

Networks

Second Edition

**Mark
Newman**

NETWORKS

Networks

Second Edition

Mark Newman

University of Michigan

OXFORD
UNIVERSITY PRESS

OXFORD

UNIVERSITY PRESS

Great Clarendon Street, Oxford, OX2 6DP,
United Kingdom

Oxford University Press is a department of the University of Oxford.
It furthers the University's objective of excellence in research, scholarship,
and education by publishing worldwide. Oxford is a registered trade mark of
Oxford University Press in the UK and in certain other countries

© Mark Newman 2018

The moral rights of the author have been asserted

First Edition published in 2010

Second Edition published in 2018

Impression: 1

All rights reserved. No part of this publication may be reproduced, stored in
a retrieval system, or transmitted, in any form or by any means, without the
prior permission in writing of Oxford University Press, or as expressly permitted
by law, by licence or under terms agreed with the appropriate reprographics
rights organization. Enquiries concerning reproduction outside the scope of the
above should be sent to the Rights Department, Oxford University Press, at the
address above

You must not circulate this work in any other form
and you must impose this same condition on any acquirer

Published in the United States of America by Oxford University Press
198 Madison Avenue, New York, NY 10016, United States of America

British Library Cataloguing in Publication Data

Data available

Library of Congress Control Number: 2018930384

ISBN 978-0-19-880509-0

Printed and bound by

CPI Group (UK) Ltd, Croydon, CR0 4YY

DOI: 10.1093/oso/9780198805090.001.0001

Links to third party websites are provided by Oxford in good faith and
for information only. Oxford disclaims any responsibility for the materials
contained in any third party website referenced in this work.

CONTENTS

Preface	ix
1 Introduction	1
I The empirical study of networks	13
2 Technological networks	14
2.1 The Internet	15
2.2 The telephone network	25
2.3 Power grids	27
2.4 Transportation networks	28
2.5 Delivery and distribution networks	29
3 Networks of information	32
3.1 The World Wide Web	32
3.2 Citation networks	37
3.3 Other information networks	41
4 Social networks	47
4.1 The empirical study of social networks	47
4.2 Interviews and questionnaires	51
4.3 Direct observation	57
4.4 Data from archival or third-party records	58
4.5 Affiliation networks	60
4.6 The small-world experiment	62
4.7 Snowball sampling, contact tracing, and random walks	65
5 Biological networks	70
5.1 Biochemical networks	70
5.2 Networks in the brain	88
5.3 Ecological networks	95

II	Fundamentals of network theory	103
6	Mathematics of networks	105
6.1	Networks and their representation	105
6.2	The adjacency matrix	106
6.3	Weighted networks	108
6.4	Directed networks	110
6.5	Hypergraphs	114
6.6	Bipartite networks	115
6.7	Multilayer and dynamic networks	118
6.8	Trees	121
6.9	Planar networks	123
6.10	Degree	126
6.11	Walks and paths	131
6.12	Components	133
6.13	Independent paths, connectivity, and cut sets	137
6.14	The graph Laplacian	142
7	Measures and metrics	158
7.1	Centrality	159
7.2	Groups of nodes	177
7.3	Transitivity and the clustering coefficient	183
7.4	Reciprocity	189
7.5	Signed edges and structural balance	190
7.6	Similarity	194
7.7	Homophily and assortative mixing	201
8	Computer algorithms	218
8.1	Software for network analysis and visualization	219
8.2	Running time and computational complexity	221
8.3	Storing network data	225
8.4	Algorithms for basic network quantities	237
8.5	Shortest paths and breadth-first search	241
8.6	Shortest paths in networks with varying edge lengths	257
8.7	Maximum flows and minimum cuts	262
9	Network statistics and measurement error	275
9.1	Types of error	276
9.2	Sources of error	278
9.3	Estimating errors	281
9.4	Correcting errors	297

10	The structure of real-world networks	304
10.1	Components	304
10.2	Shortest paths and the small-world effect	310
10.3	Degree distributions	313
10.4	Power laws and scale-free networks	317
10.5	Distributions of other centrality measures	330
10.6	Clustering coefficients	332
10.7	Assortative mixing	335
III	Network models	341
11	Random graphs	342
11.1	Random graphs	343
11.2	Mean number of edges and mean degree	345
11.3	Degree distribution	346
11.4	Clustering coefficient	347
11.5	Giant component	347
11.6	Small components	355
11.7	Path lengths	360
11.8	Problems with the random graph	364
12	The configuration model	369
12.1	The configuration model	370
12.2	Excess degree distribution	377
12.3	Clustering coefficient	381
12.4	Locally tree-like networks	382
12.5	Number of second neighbors of a node	383
12.6	Giant component	384
12.7	Small components	391
12.8	Networks with power-law degree distributions	395
12.9	Diameter	399
12.10	Generating function methods	401
12.11	Other random graph models	416
13	Models of network formation	434
13.1	Preferential attachment	435
13.2	The model of Barabási and Albert	448
13.3	Time evolution of the network and the first mover effect	451
13.4	Extensions of preferential attachment models	458
13.5	Node copying models	472
13.6	Network optimization models	479

IV Applications	493
14 Community structure	494
14.1 Dividing networks into groups	495
14.2 Modularity maximization	498
14.3 Methods based on information theory	515
14.4 Methods based on statistical inference	520
14.5 Other algorithms for community detection	529
14.6 Measuring algorithm performance	538
14.7 Detecting other kinds of network structure	551
15 Percolation and network resilience	569
15.1 Percolation	569
15.2 Uniform random removal of nodes	571
15.3 Non-uniform removal of nodes	586
15.4 Percolation in real-world networks	593
15.5 Computer algorithms for percolation	594
16 Epidemics on networks	607
16.1 Models of the spread of infection	608
16.2 Epidemic models on networks	624
16.3 Outbreak sizes and percolation	625
16.4 Time-dependent properties of epidemics on networks	645
16.5 Time-dependent properties of the SI model	646
16.6 Time-dependent properties of the SIR model	660
16.7 Time-dependent properties of the SIS model	667
17 Dynamical systems on networks	675
17.1 Dynamical systems	676
17.2 Dynamics on networks	685
17.3 Dynamics with more than one variable per node	694
17.4 Spectra of networks	698
17.5 Synchronization	701
18 Network search	710
18.1 Web search	710
18.2 Searching distributed databases	713
18.3 Sending messages	718
References	732
Index	751

PREFACE

The scientific study of networks, such as computer networks, biological networks, and social networks, is an interdisciplinary field that combines ideas from mathematics, physics, biology, computer science, statistics, the social sciences, and many other areas. The field has benefited enormously from the wide range of viewpoints brought to it by practitioners from so many different disciplines, but it has also suffered because human knowledge about networks is dispersed across the scientific community and researchers in one area often do not have ready access to discoveries made in another. The goal of this book is to bring our knowledge of networks together and present it in consistent language and notation, so that it becomes a coherent whole whose elements complement one another and in combination teach us more than any single element can alone.

The book is divided into four parts. Following a short introductory chapter, Part I describes the basic types of networks studied by present-day science and the empirical techniques used to determine their structure. Part II introduces the fundamental tools used in the study of networks, including the mathematical methods used to represent network structure, measures and statistics for quantifying network structure, and computer algorithms for calculating those measures and statistics. Part III describes mathematical models of network structure that can help us predict the behavior of networked systems and understand their formation and growth. And Part IV describes applications of network theory, including models of network resilience, epidemics taking place on networks, and network search processes.

The technical level of the presentation varies among the parts, Part I requiring virtually no mathematical knowledge for its comprehension, while Part II requires a grasp of linear algebra and calculus at the undergraduate level. Parts III and IV are mathematically more advanced and suitable for advanced undergraduates, postgraduates, and researchers working in the field. The book could thus be used as the basis of a taught course at various levels. A less technical course suitable for those with moderate mathematical knowledge might cover the material of Chapters 1 to 10, while a more technical course for

advanced students might cover the material of Chapters 6 to 13 and selected material thereafter. Each chapter from Part II onwards is accompanied by a selection of exercises that can be used to test the reader's understanding of the material.

The study of networks is a rapidly advancing field and this second edition of the book includes a significant amount of new material, including sections on multilayer networks, network statistics, community detection, complex contagion, and network synchronization. The entire book has been thoroughly updated to reflect recent developments in the field and many new exercises have been added throughout.

Over its two editions this book has been some years in the making and many people have helped me with it during that time. I must thank my ever-patient editor Sonke Adlung, with whom I have worked on various book projects for more than 25 years, and whose constant encouragement and wise advice have made working with him and Oxford University Press a real pleasure. Thanks are also due to Melanie Johnstone, Viki Kapur, Charles Lauder, Alison Lees, Emma Lonie, April Warman, and Ania Wronski for their help with the final stages of bringing the book to print.

I have benefited greatly during the writing of the book from the conversation, comments, suggestions, and encouragement of many colleagues and friends. They are, sadly, too numerous to mention exhaustively, but special thanks must go to Edoardo Airoidi, Robert Axelrod, Steve Borgatti, Elizabeth Bruch, Duncan Callaway, François Caron, Aaron Clauset, Robert Deegan, Jennifer Dunne, Betsy Foxman, Linton Freeman, Michelle Girvan, Mark Handcock, Petter Holme, Jon Kleinberg, Alden Klovdahl, Liza Levina, Lauren Meyers, Cris Moore, Lou Pecora, Mason Porter, Sidney Redner, Gesine Reinert, Martin Rosvall, Cosma Shalizi, Steve Strogatz, Duncan Watts, Doug White, Lenka Zdeborová, and Bob Ziff, as well as to the many students and other readers whose feedback helped iron out a lot of rough spots, particularly Michelle Adan, Alejandro Balbin, Ken Brown, George Cantwell, Judson Caskey, Rachel Chen, Chris Fink, Massimo Franceschet, Milton Friesen, Michael Gastner, Martin Gould, Timothy Griffin, Ruthi Hortsch, Shi Xiang Lam, Xiaoning Qian, Harry Richman, Puck Rombach, Tyler Rush, Snehal Shekatkar, Weijing Tang, Robb Thomas, Jane Wang, Paul Wellin, Daniel Wilcox, Yongsoo Yang, and Dong Zhou. I would also especially like to thank Brian Karrer, who read the entire book in draft form and gave me many pages of thoughtful and thought-provoking comments, as well as spotting a number of mistakes and typos. Responsibility for any remaining mistakes in the book of course rests entirely

with myself, and I welcome corrections from readers.

Finally, my heartfelt thanks go to my wife Carrie for her continual encouragement and support during the writing of this book. Without her the book would still have been written but I would have smiled a lot less.

Mark Newman
Ann Arbor, Michigan
June 12, 2018

CHAPTER 1

INTRODUCTION

*A short introduction to networks
and why we study them*

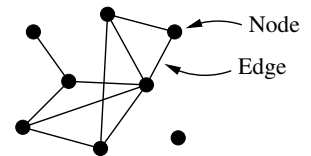
A NETWORK is, in its simplest form, a collection of points joined together in pairs by lines. In the nomenclature of the field a point is referred to as a *node* or *vertex*¹ and a line is referred to as an *edge*. Many systems of interest in the physical, biological, and social sciences can be thought of as networks and, as this book aims to show, thinking of them in this way can lead to new and useful insights.

We begin in this first chapter with a brief introduction to some of the most commonly studied types of networks and their properties. All the topics in this chapter are covered in greater depth later in the book.

EXAMPLES OF NETWORKS

Networks of one kind or another crop up in almost every branch of science and technology. We will encounter a huge array of interesting examples in this book. Purely for organizational purposes, we will divide them into four broad categories: technological networks, information networks, social networks, and biological networks.

A good example of a technological network is the Internet, the computer data network in which the nodes are computers and the edges are data connections between them, such as optical fiber cables or telephone lines. Figure 1.1



A small network composed of eight nodes and ten edges.

¹Plural: vertices.

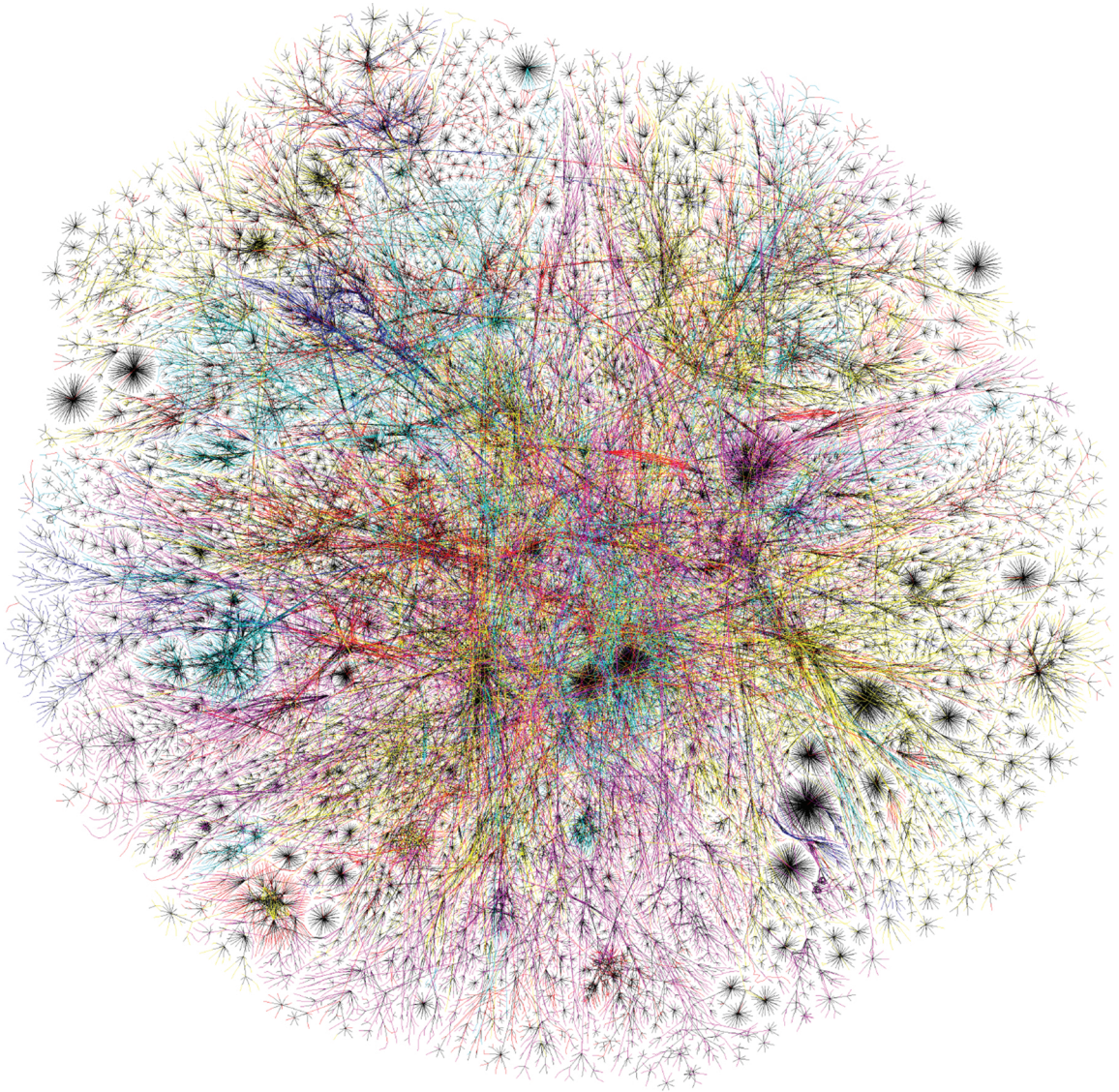


Figure 1.1: The network structure of the Internet. The nodes in this representation of the Internet are “class C subnets”—groups of computers with similar Internet addresses that are usually under the management of a single organization—and the connections between them represent the routes taken by Internet data packets as they hop between subnets. The geometric positions of the nodes in the picture have no special meaning; they are chosen simply to give a pleasing layout and are not related, for instance, to geographic position of the nodes. The structure of the Internet is discussed in detail in Section 2.1. Figure created by the Opte Project (<http://www.opte.org>). Reproduced with permission.

shows a picture of the structure of the Internet, a snapshot of the network as it was in 2003, reconstructed by observing the paths taken across the network by a large number of Internet data packets. It is a curious fact that although the Internet is a man-made and carefully engineered network, we don't know exactly what its structure is because it was built by many different groups of people with only limited knowledge of each other's actions and little centralized control. Our best current data on its structure are therefore derived from experimental measurements, such as those that produced this figure, rather than from any centrally held map or repository of knowledge.

There are a number of practical reasons why we might want to study the network structure of the Internet. The function of the Internet is to transport data between computers (and other devices) in different parts of the world, which it does by dividing the data into separate packets and shipping them from node to node across the network until they reach their intended destination. The network structure of the Internet will affect how efficiently it performs this function, and if we know that structure we can address many questions of practical relevance. How should we choose the route by which data are transported? Is the shortest route, geographically speaking, always necessarily the fastest? If not, then what is, and how can we find it? How can we avoid bottlenecks in the traffic flow that might slow things down? What happens when a node or an edge fails (which they do with some regularity)? How can we devise schemes to route around such failures? If we have the opportunity to add new capacity to the network, where should it be added?

Other examples of technological networks include the telephone network, networks of roads, rail lines, or airline routes, and distribution networks such as the electricity grid, water lines, oil or gas pipelines, or sewerage pipes. Each of these networks raises questions of their own: what is their structure, how does it affect the function of the system, and how can we design or change the structure to optimize performance? In some cases, such as airline routes, networks are already highly optimized; in others, such as the road network, the structure may be largely a historical accident and is in some cases far from optimal.

Our second class of networks are the information networks, a more abstract class that represents the network structure of bodies of information. The classic example is the World Wide Web. We discussed the Internet above, but the Web is not the same thing as the Internet, even though the two words are often used interchangeably in casual speech. The Internet is a physical network of computers linked by actual cables (or sometimes radio links) running between them. The Web, on the other hand, is a network of web pages and the links between them. The nodes of the World Wide Web are the web pages and the

We look at the Internet in more detail in Section 2.1.

Information networks are discussed at length in Chapter 3.

The World Wide Web is discussed in more detail in Section 3.1.

edges are “hyperlinks,” the highlighted snippets of text or push-buttons on web pages that we click on to navigate from one page to another. A hyperlink is purely a software construct; you can link from your web page to a page that lives on a computer on the other side of the world just as easily as you can to a friend down the hall. There is no physical structure, like an optical fiber, that needs to be built when you make a new link. The link is merely an address that tells the computer where to look next when you click on it. Thus the network structure of the Web and the Internet are completely distinct.

Abstract though it may be, the World Wide Web, with its billions of pages and links, has proved enormously useful, not to mention profitable, and the structure of the network is of substantial interest. Since people tend to add hyperlinks between pages with related content, the link structure of the Web reveals something about relationships between content and topics. Arguably, the structure of the Web could be said to reflect the structure of human knowledge. What’s more, people tend to link more often to pages they find useful than to those they do not, so that the number of links pointing to a page can be used as a measure of its usefulness. A more sophisticated version of this idea lies behind the operation of the popular web search engine *Google*, as well as some others.

The mechanics of web search are discussed in Section 18.1.

The Web also illustrates another concept of network theory, the *directed network*. Hyperlinks on the Web run in one specific direction, from one web page to another. You may be able to click a link on page A and get to page B, but there is no requirement that B has a link back to A again. (It might contain such a link but it doesn’t have to.) One says that the edges in the World Wide Web are *directed*, running from the linking page to the linked.

Another much-studied example of an information network is a citation network, such as the network of citations between academic journal articles. Academic articles typically include a bibliography of references to other previously published articles, and one can think of these references as forming a network in which the articles are the nodes and there is a directed edge from article A to article B if A cites B in its bibliography. As with the World Wide Web, one can argue that such a network reflects, at least partially, the structure of the body of knowledge contained in the articles, with citations between articles presumably indicating related content. Indeed there are many similarities between the Web and citation networks and a number of the techniques developed for understanding and searching the Web have in recent years started to be applied to citation networks too, to help scientists and others filter the vast amount of published research and data to find useful papers.

Our third broad class of networks are the social networks. When one talks about “social networks” today, most of us think of online services such as

Facebook or *Twitter*, but in the scientific literature the term is used much more broadly to encompass any network in which the nodes are people (or sometimes groups of people, such as firms or teams) and the edges between them are social connections of some kind, such as friendship, communication, or collaboration. The field of sociology has perhaps the longest and best developed tradition of the empirical study of networks as they occur in the real world, and many of the mathematical and statistical tools used in the study of networks are borrowed, directly or indirectly, from sociologists.

Figure 1.2 shows a famous example of a social network from the sociology literature, Wayne Zachary’s “karate club” network. This network represents the pattern of friendships among the members of a karate club at a North American university, reconstructed from observations of social interactions between them. Sociologists have performed a huge number of similar studies over the decades, including studies of friendship patterns among CEOs of corporations, doctors, monks, students, and conference participants, and networks of who works with whom, who does business with whom, who seeks advice from whom, who socializes with whom, and who sleeps with whom. Such studies, in which data are typically collected by hand, are quite arduous, so the networks they produce are usually small, like the one in Fig. 1.2, which has just 34 nodes. But in recent years, much larger social networks have been assembled using, for instance, online data from Facebook and similar services. At the time of writing, Facebook had over two billion users worldwide—more than a quarter of the population of the world—and information on the connection patterns between all of them. Many online social networking companies, including Facebook, have research divisions that collaborate with the academic community to do research on social networks using their vast data resources.

Our fourth and final class of networks is biological networks. Networks occur in range of different settings in biology. Some are physical networks like neural networks—the connections between neurons in the brain—while others are more abstract. In Fig. 1.3 we show a picture of a “food web,” an ecological network in which the nodes are species in an ecosystem and the edges represent predator–prey relationships between them. That is, a pair of species is connected by an edge in this network if one species eats the other. The study of food webs can help us understand and quantify many ecological phenomena, particularly concerning energy and carbon flows and the interdependencies

Social networks are discussed in more depth in Chapter 4.

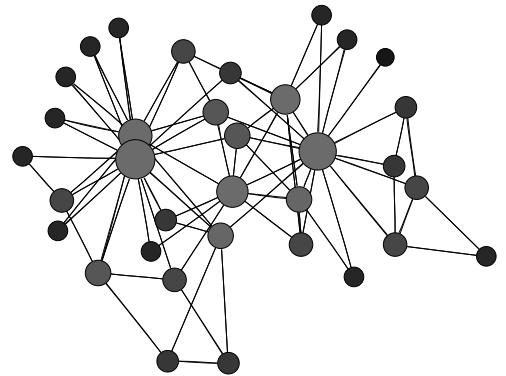


Figure 1.2: Friendship network between members of a club. This social network from a study conducted in the 1970s shows the pattern of friendships between the members of a karate club at an American university. The data were collected and published by Zachary [479].

Neural networks are discussed in Section 5.2 and food webs in Section 5.3.

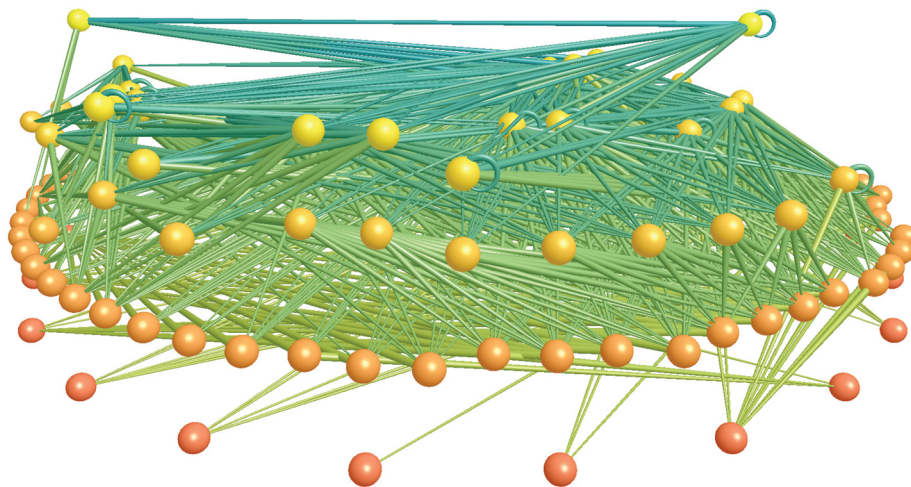


Figure 1.3: The food web of Little Rock Lake, Wisconsin. This elegant picture summarizes the known predatory interactions between species in a freshwater lake in the northern United States. The nodes represent the species and the edges run between predator-prey species pairs. The vertical position of the nodes represents, roughly speaking, the trophic level of the corresponding species. The figure was created by Richard Williams and Neo Martinez [321].

between species. Food webs also provide us with another example of a directed network, like the World Wide Web and citation networks discussed previously. If species A eats species B then probably B does not also eat A, so the relationship between the two is a directed one.

Biochemical networks are discussed in detail in Section 5.1.

An example metabolic network map appears as Fig. 5.2 on page 75.

Another class of biological networks are the biochemical networks. These include metabolic networks, protein-protein interaction networks, and genetic regulatory networks. A metabolic network, for instance, is a representation of the pattern of chemical reactions that fuel the cells in an organism. The reader may have seen the wallcharts of metabolic reactions that adorn the offices of some biochemists, incredibly detailed maps with hundreds of tiny inscriptions linked by a maze of arrows. The inscriptions—the nodes in this network—are metabolites, the chemicals involved in metabolism, and the arrows—directed edges—are reactions that turn one metabolite into another. The representation of reactions as a network is one of the first steps towards making sense of the bewildering array of biochemical data generated by current experiments in biochemistry and molecular genetics.

These are just a few examples of the types of networks that will concern us

in this book. These and many others are studied in more detail in the following chapters.

WHAT CAN WE LEARN FROM NETWORKS?

Networks capture the pattern of interactions between the parts of a system. It should come as no surprise (although in some fields it is a relatively recent realization) that the pattern of interactions can have a big effect on the behavior of a system. The pattern of connections between computers on the Internet, for instance, affects the routes that data take over the network and hence the efficiency with which the network transports those data. The connections in a friendship network affect how people learn, form opinions, and gather news, as well as other less obvious phenomena, such as the spread of disease. Unless we know something about the structure of these networks, we cannot hope to understand fully how the corresponding systems work.

A network is a simplified representation that reduces a system to an abstract structure or *topology*, capturing only the basics of connection patterns and little else. The systems studied can, and often do, have many other interesting features not represented by the network—the detailed behaviors of individual nodes, such as computers or people, for instance, or the precise nature of the interactions between them. Some of these subtleties can be captured by embroidering the network with labels on the nodes or edges, such as names or strengths of interactions, but even so a lot of information is usually lost in the process of reducing a full system to a network representation. This has some disadvantages but it has advantages as well.

Scientists in a wide variety of fields have, over the years, developed an extensive set of mathematical and computational tools for analyzing, modeling, and understanding networks. Some of these tools start from a simple network topology—a set of nodes and edges—and after some calculation tell you something potentially useful about the network: which is the best connected node, say, or how similar two nodes are to one another. Other tools take the form of network models that can make mathematical predictions about processes taking place on networks, such as the way traffic will flow over the Internet or the way a disease will spread through a community. Because they work with networks in their abstract form, tools such as these can be applied to almost any system that has a network representation. Thus, if there is a system you are interested in, and it can usefully be represented as a network, then there are hundreds of ready-made tools out there, already fully developed and well understood, that you can immediately apply to your system. Not all of them will necessarily give useful results—which measurements or calculations are

Some common network extensions and variants are discussed in Chapter 6.

useful for a particular system depends on what the system is and does and on what specific questions you are trying to answer about it. Still, if you have a well-posed question about a networked system there will, in many cases, already be a tool available that will help you address it.

Networks are thus a general means for representing the structure of a system that creates a bridge between empirical data and a large toolkit of powerful analysis techniques. In this book we discuss many examples of specific networks in different fields, along with techniques for their analysis drawn from mathematics, physics, the computer and information sciences, the social sciences, biology, and elsewhere. In doing so, we will bring together a wide range of ideas and expertise from many disciplines to build a comprehensive understanding of the science of networks.

PROPERTIES OF NETWORKS

Perhaps the most fundamental question we can ask about networks is this: if we know the shape of a network, what can we learn about the nature and function of the system it describes? In other words, how are the structural features of a network related to the practical issues we care about? This question is essentially the topic of this entire book, and we are not going to answer it in this chapter alone. Let us, however, look briefly here at a few representative concepts, to get a feel for the kinds of ideas we will be dealing with.

A first step in analyzing the structure of a network is often to make a picture of it. Figures 1.1, 1.2, and 1.3 are typical examples. Each of them was generated by a specialized computer program designed for network visualization and there are many such programs available, both commercially and for free, if you want to produce pictures like these for yourself. Visualization can be an extraordinarily useful tool in the analysis of network data, allowing one to instantly see important structural features that would otherwise be difficult to pick out of the raw data. The human eye is enormously gifted at discerning patterns, and visualizations allow us to put this gift to work on our network problems.

On the other hand, direct visualization of networks is only really useful for networks up to a few hundreds or thousands of nodes, and for networks that are relatively sparse, meaning that the number of edges is quite small. If there are too many nodes or edges then pictures of the network will be too complicated for the eye to comprehend and their usefulness becomes limited. Many of the networks that scientists are interested in today have hundreds of thousands or even millions of nodes, which means that visualization is not of much help and we need to employ other techniques to understand them. Moreover, while the

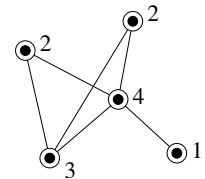
eye is definitely a powerful tool for data analysis, it is not a wholly reliable one, sometimes failing to pick out important patterns in data or even seeing patterns where they don't exist. To address these issues, network theory has developed a large toolchest of measures and metrics that can help us understand what networks are telling us, even in cases where useful visualization is impossible or unreliable.

An example of a useful (and widely used) class of network metrics are the *centrality* measures. Centrality quantifies how important nodes are in a network, and social network analysts in particular have expended considerable effort studying it. There are, of course, many different possible concepts or definitions of what it means for a node to be central in a network, and there are correspondingly many centrality measures. Perhaps the simplest of them is the measure called *degree*. The degree of a node in a network is the number of edges attached to it. In a social network of friendships, for instance, such as the network of Fig. 1.2, the degree of an individual is the number of friends he or she has within the network. For the Internet degree would be the number of data connections a computer has. In many cases the nodes with the highest degrees in a network, those with the most connections, also play major roles in the functioning of the system, and hence degree can be a useful guide for focusing our attention on the system's most important elements.

In undirected networks degree is just a single number, but in directed networks nodes have two different degrees, *in-degree* and *out-degree*, corresponding to the number of edges pointing inward and outward respectively. For example, the in-degree of a web page is the number of other pages that link to it, while the out-degree is the number of pages to which it links. We have already mentioned one example of how centrality can be put to use on the Web to answer an important practical question: by counting the number of links a web page gets—the in-degree of the page—a search engine can identify pages that are likely to contain useful information.

A further observation concerning degree is that many networks are found to contain a small but significant number of “hubs”—nodes of unusually high degree. Social networks, for instance, often contain a few individuals with an unusually large number of acquaintances. The Web has a small fraction of websites with a very large number of links. There are a few metabolites that take part in a very large number of metabolic processes. A major topic of research in recent years has been the investigation of the effects of hubs on the performance and behavior of networked systems. A wide range of results, both empirical and theoretical, indicate that hubs can have a disproportionate effect, particularly on network resilience and transport phenomena, despite being few in number.

See Chapter 7 for further discussion of centrality measures.



The number beside each node in this small network indicates the node's degree.

Hubs are discussed further in Section 10.3.

The small-world effect is discussed further in Sections 4.6 and 10.2.

Another example of a network concept that arises repeatedly and has real practical implications is the so-called *small-world effect*. Given a network, one can ask what the shortest distance is, through the network, between a given pair of nodes. In other words, what is the minimum number of edges one would have to traverse in order to get from one node to the other? For instance, your immediate friend would have distance 1 from you in a network of friendships, while a friend of a friend would have distance 2. It has been found empirically (and can be proven mathematically in some cases) that the mean distance between node pairs in many networks is very short, often no more than a dozen steps or so, even for networks with millions of nodes or more. Although first studied in the context of friendship networks, this small-world effect appears to be widespread, occurring in essentially all types of networks. In popular culture it is referred to as the “six degrees of separation,” after a successful stage play and film of the same name in which the effect is discussed. The (semi-mythological) claim is that you can get from anyone in the world to anyone else via a sequence of no more than five intermediate acquaintances—six steps in all.

The small-world effect has substantial repercussions. For example, news and gossip spread over social networks—if you hear an interesting rumor from a friend, you may pass it on to your other friends, and they in turn may pass it on to theirs, and so forth. Clearly a rumor will spread faster and further if it only takes six steps to reach anyone in the world than if it takes a hundred, or a million. And indeed it is a matter of common experience that a suitably scandalous rumor can reach the ears of an entire community in what seems like the blink of an eye.

Or consider the Internet. One of the reasons the Internet functions at all is because any computer on the network is only a few hops across the network from any other. Typical routes taken by data packets over the Internet rarely have more than about twenty hops, and certainly the performance of the network would be much worse if packets had to make a thousand hops instead. In effect, our ability to receive data near instantaneously from anywhere in the world is a direct consequence of the small-world effect.

Community structure in networks is discussed in detail in Chapter 14.

A third example of a network phenomenon of practical importance is the occurrence of clusters or communities in networks. We are most of us familiar with the idea that social networks break up into subcommunities. In friendship networks, for instance, one commonly observes groups of close friends within the larger, looser network of passing acquaintances. Similar clusters occur in other types of network as well. The Web contains clusters of web pages that all link to one another, perhaps because they are about the same topic, or they all belong to the same company. Metabolic networks contain groups of metabolites

that interact with one another to perform certain biochemical tasks. And if it is the case that clusters or groups correspond to functional divisions in this way, then we may be able to learn something by taking a network and decomposing it into its constituent clusters. The way a network breaks apart can reveal levels and concepts of organization that are not easy to see by other means.

The detection and analysis of clusters in networks is an active topic at the frontier of current networks research, holding promise for exciting applications in the future.

OUTLINE OF THIS BOOK

This book is divided into four parts. In the first part, consisting of Chapters 2 to 5, we introduce the various types of network encountered in the real world, including technological, social, and biological networks, and the empirical techniques used to discover their structure. Although it is not the purpose of this book to describe any one particular network in great detail, the study of networks is nonetheless firmly founded on empirical observations and a good understanding of what data are available and how they are obtained is immensely helpful in understanding the science of networks as it is practiced today.

The second part of the book, Chapters 6 to 10, introduces the fundamental theoretical ideas and methods on which our current understanding of networks is based. Chapter 6 describes the basic mathematics used to capture network ideas, while Chapter 7 describes the measures and metrics we use to quantify network structure. Chapter 8 describes the computer methods that are crucial to practical calculations on today's large networks, Chapter 9 describes methods of network statistics and the role of errors and uncertainty in network studies, and Chapter 10 describes some of the intriguing patterns and principles that emerge when we apply all of these ideas to real-world network data.

In the third part of the book, Chapters 11 to 13, we look at mathematical models of networks, including both traditional models, such as random graphs and their extensions, and newer models, such as models of growing networks and community structure. The material in these chapters forms a central part of the canon of the field and has been the subject of a vast amount of published scientific research.

Finally, in the fourth and last part of the book, Chapters 14 to 18, we look at applications of network theory to a range of practical questions, including community detection, network epidemiology, dynamical systems, and network search processes. Research is less far advanced on these topics than it is in other areas of network science and there is much we do not know. The final chapters

INTRODUCTION

of the book probably raise at least as many questions as they answer, but this, surely, is a good thing. For those who would like to get involved, there are plenty of fascinating open problems waiting to be addressed.

PART I

THE EMPIRICAL STUDY OF
NETWORKS

CHAPTER 2

TECHNOLOGICAL NETWORKS

A discussion of engineered networks like the Internet and the power grid and methods for determining their structure

The four classes are not rigorously defined and there is, as we will see, some overlap between them, with some networks plausibly belonging to two or more classes. Nonetheless, the division into classes is a useful one, since networks in the same class are often treated using similar techniques or ideas.

IN THE next four chapters we describe and discuss some of the most commonly studied networks, dividing them into four broad classes—technological networks, information networks, social networks, and biological networks. For each class we list some important examples and examine the techniques used to measure their structure.

It is not our intention in this book to study any one network in great detail. Plenty of other books do that. Nonetheless, network science is concerned with understanding and modeling the behavior of real-world systems and observational data are the starting point for essentially all the developments of the field, so it will be useful to have a grasp of the types of networks commonly studied and the data that describe them. In this chapter we look at technological networks, the physical infrastructure networks that form the backbone of modern technological societies. Perhaps the most celebrated such network—and a relatively recent entry in the field—is the Internet, the global network of data connections that links computers and other information systems together. Section 2.1 is devoted to a discussion of the Internet. A number of other important examples of technological networks, including power grids, transportation networks, delivery and distribution networks, and telephone networks, are discussed in subsequent sections.

Networks, 2nd edition. Mark Newman, Oxford University Press (2018). © Mark Newman.
DOI: 10.1093/oso/9780198805090.001.0001

2.1 THE INTERNET

The Internet is the worldwide network of physical data connections between computers, phones, tablets, and other devices. The Internet is a *packet-switched* data network, meaning that messages sent over it are broken up into *packets*, small chunks of data, that are sent separately over the network and reassembled into a complete message again at the other end. The format of the packets follows a standard known as the *Internet Protocol* (IP) and includes an *IP address* in each packet that specifies the packet's destination, so that it can be routed correctly across the network.

The simplest network representation of the Internet (there are others, which we will discuss shortly) is one in which the nodes of the network represent computers and other devices, and the edges represent data connections between them, such as optical fiber lines or wireless connections. In fact, ordinary computers and other consumer devices mostly occupy the nodes on the “outside” of the network, the end points (or starting points) of data flows, and do not act as intermediate points between others. (Indeed, most end-user devices only have a single connection to the Net, so it would not be possible for them to lie on a path between any others.) The “interior” nodes of the Internet are primarily *routers*, powerful special-purpose machines at the junctions between data lines that receive data packets and forward them in one direction or another towards their intended destination (essentially larger versions of the network router you might have in your home).

The general overall shape of the Internet is shown, in schematic form, in Fig. 2.1. The network is composed of three levels or circles of nodes. The innermost circle, the core of the network, is called the *backbone* and contains the trunk lines that provide long-distance high-bandwidth data transport across the globe, along with the high-performance routers and switching centers that link the trunk lines together. The trunk lines are the highways of the Internet, built with the fastest fiber optic connections available (and improving all the time). The backbone is owned and operated by a set of *network backbone providers* (NBPs), who are primarily national governments and major telecommunications companies such as Level 3 Communications, Cogent, NTT, and others.

The second circle of the Internet is composed of *Internet service providers* or ISPs—commercial companies, governments, universities, and others who contract with NBPs for connection to the backbone and then resell or otherwise provide that connection to end users, the ultimate consumers of Internet bandwidth, who form the third circle—businesses, government offices, academic institutions, people in their homes, and so forth. As Fig. 2.1 shows, the ISPs

The Internet should not be confused with the World Wide Web, a virtual network of web pages and hyperlinks, which we discuss separately in Section 3.1.

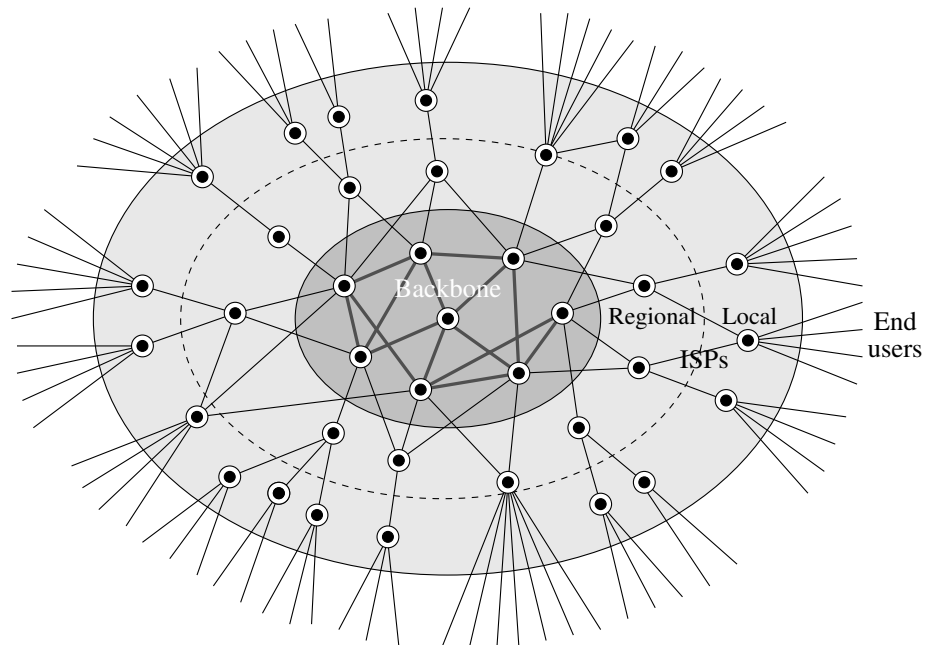


Figure 2.1: A schematic depiction of the structure of the Internet. The nodes and edges of the Internet fall into a number of different classes: the backbone of high-bandwidth long-distance connections; the ISPs, who connect to the backbone and who are divided roughly into regional (larger) and local (smaller) ISPs; and the end users—home users, companies, and so forth—who connect to the ISPs.

are further subdivided into *regional ISPs* and *local* or *consumer ISPs*, the former being larger organizations whose primary customers are the local ISPs, who in turn sell network connections to the end users. This distinction is somewhat blurred however, because large consumer ISPs, such as AT&T or British Telecom, often act as their own regional ISPs (and some may be backbone providers as well).

The network structure of the Internet is not dictated by any central authority. Protocols and guidelines are developed by an informal volunteer organization called the Internet Engineering Task Force, but one does not have to apply to any central Internet authority for permission to build a new spur on the network, or to take one out of service.

One of the remarkable features of the Internet is the scheme used for routing data across the network, in which the paths that packets take are determined by automated negotiation among Internet routers under a system called the

Border Gateway Protocol (BGP). BGP is designed in such a way that if new nodes or edges are added to the network or old ones disappear, either permanently or temporarily, routers will take note and adjust their routing policy appropriately. There is a certain amount of human oversight involved, to make sure the system keeps running smoothly, but no “Internet government” is needed to steer things from on high; the system organizes itself by the combined actions of many local and essentially autonomous computer systems.

While this is an excellent feature of the system from the point of view of robustness and flexibility, it is a problem for those who want to study the structure of the Internet. If there were a central Internet government with a complete map of the system, then the job of determining the network structure would be easy—one would just look at the map. But there is no such organization and no such map. Instead the network’s structure must be determined by experimental measurements. There are two primary methods for doing this. The first uses “traceroute”; the second uses BGP.

2.1.1 MEASURING INTERNET STRUCTURE USING TRACEROUTE

There is currently no simple means by which to probe the network structure of the Internet directly. We can, however, quite easily discover the particular path taken by data packets sent from one computer to another on the Internet. The standard tool for doing this is called *traceroute*.

Each Internet data packet contains, among other things, a destination address, which says where it is going; a source address, which says where it started from; and a *time-to-live* (TTL). The TTL is a number that specifies the maximum number of hops that the packet can make to get to its destination, a hop being the traversal of one edge in the network. At every hop, the TTL is decreased by one, and if it reaches zero the packet is discarded, meaning it is deleted and not forwarded any further over the network. A message is also then transmitted back to the sender informing them that the packet was discarded and where it got to. In this way the sender is alerted if data is lost, allowing them to resend the contents of the packet if necessary. The TTL exists mainly as a safeguard to prevent packets from losing their way on the Internet and wandering around forever, but we can make use of it to track packet progress as well. The idea is as follows.

First, we send out a packet with the destination address of the network node we are interested in and a TTL of 1. The packet makes a single hop to the first router along the way, its TTL is decreased to 0, the packet is discarded by the router, and a message is returned to us telling us, among other things, the IP address of the router. We record this address and then repeat the process

with a TTL of 2. This time the packet makes two hops before dying and the returned message tells us the IP address of the second router along the path. The process is repeated with larger and larger TTL until the destination is reached, and the set of IP addresses received as a result tells us the entire route taken to get there.¹ There are standard software tools that will perform the complete procedure automatically and print out the list of IP addresses for us. On many operating systems the tool that does this is called “traceroute.”²

We can use traceroute (or a similar tool) to probe the network structure of the Internet. The idea is to assemble a large data set of traceroute paths between many different pairs of points on the Internet. With luck, most of the edges in the network (though usually not all of them) will appear in at least one of these paths, and the combination of all of them together should give a reasonably complete picture of the network. Early studies, for the sake of expediency, limited themselves to paths starting from just a few source computers, but more recent ones make use of distributed collections of thousands of sources to develop a very complete picture of the network.

The paths from any single source to a set of destinations form a branching structure as shown schematically in Figs. 2.2a, b, and c.³ The source computers should, ideally, be well distributed over the network. If they are close together then there may be substantial overlap between the paths to distant nodes, meaning that they will needlessly duplicate each other’s efforts rather than returning independent measurements.

Once one has a suitable set of traceroute paths, a simple union of them gives us our snapshot of the network structure—see Fig. 2.2d. That is, we create a node in our network for every unique IP address that appears at least once in any of the paths and an edge between any pair of addresses that fall on adjacent steps of any path. As hinted above, it is unlikely that this procedure

¹We are assuming that each packet takes the same route to the destination. It is possible, though rare, for different packets to take different routes, in which case the set of IP addresses returned by the traceroute procedure will not give a correct path through the network. This can happen, for instance, if congestion patterns along the route vary significantly while the procedure is being performed, causing the network to reroute packets along less congested paths. Serious Internet mapping experiments perform repeated traceroute measurements to minimize the errors introduced by effects such as these.

²On the Windows operating system it is called “tracert.” On some Linux systems it is called “tracepath.”

³If there were a unique best path to every node, then the set of paths would be a “tree,” meaning it would contain no loops. (See Section 6.8 for a discussion of trees.) Because of the way routing algorithms work, however, this is not always the case in practice—two routes that originate at the same point and pass through the same node on the way to their final destination can still take different routes to get to that node, so that the set of paths can contain loops.

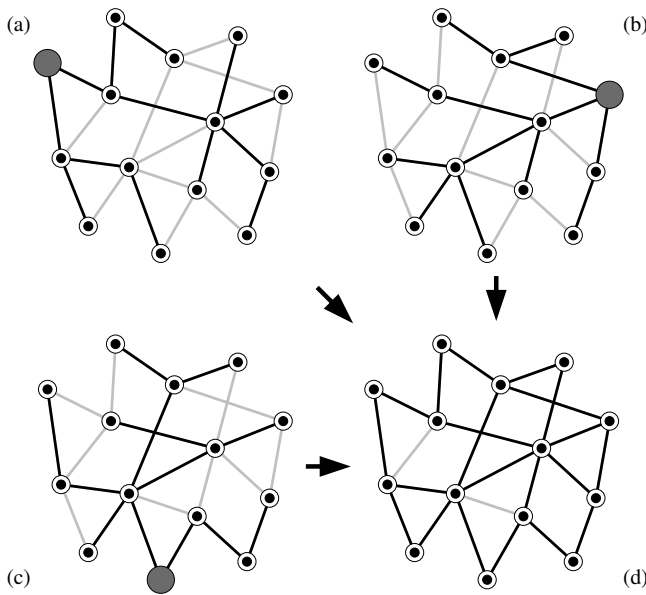


Figure 2.2: Reconstruction of the topology of the Internet from traceroute data. In (a), (b), and (c) we show in bold the edges that fall along traceroute paths starting from the three highlighted source nodes. In (d) we form the union of these edges to make a picture of the overall network topology. Note that a few edges are still missing from the final picture (the remaining gray edges in (d)) because they happen not to appear in any of the three individual traceroute data sets.

will find all the edges in the network (see Fig. 2.2d again), and for studies based on small numbers of sources there can be significant biases in the sampling of edges [3,284]. However, better and better data sets are becoming available as time passes, and it is believed that we now have a reasonably complete picture of the shape of the Internet.

In fact, complete (or near-complete) representations of the Internet of the kind described here can be cumbersome to work with and are typically not used directly for network studies. There are billions of distinct IP addresses in use on the Internet at any one time, with many of those corresponding to end-user devices that appear or disappear as the devices are turned on or off or connections to the Internet are made or broken. Most studies of the Internet ignore end users and restrict themselves to just the routers, in effect concentrating on the inner zones in Fig. 2.1 and ignoring the outermost one. We will refer to such maps of the Internet as representations at the *router level*. The nodes in the network are routers, and the edges between them are network connections.

It may appear strange to ignore end-user devices, since the end users are, after all, the entire reason for the Internet's existence in the first place. However, it is the structure of the network at the router level that is responsible for most aspects of the performance, robustness, and efficiency of the network, that dictates the patterns of traffic flow on the network, and that forms the focus of most work on Internet structure and design. To the extent that these are the issues of scientific interest, therefore, it makes sense to concentrate our efforts on the router-level structure.

An example of a study of the topology of the Internet at the router level is that of Faloutsos *et al.* [168], who looked at the "degree distribution" of the network and discovered it to follow, approximately, a power law. We discuss degree distributions and power laws in networks in more detail in Section 10.4.

Even after removing all or most end users from the network, the structure of the Internet at the router level may still be too detailed for our purposes. Often we would like a more coarse-grained representation of the network that gives us a broader overall picture of network structure. Such representations can be created by grouping sets of IP addresses together into single nodes. Three different ways of grouping addresses are in common use, giving rise to three different coarse-grained representations, at the level of subnets, domains, and autonomous systems.

A *subnet* is a group of IP addresses defined as follows. IP addresses consist of four numbers, each one in the range from 0 to 255 (eight bits in binary) and typically written in a string separated by periods or dots.⁴ For example, the IP address of the main web server at the author's home institution, the University of Michigan, is 141.211.243.44. IP addresses are allocated to organizations in blocks. The University of Michigan, for instance, owns (among others) all the addresses of the form 141.211.243.xxx, where "xxx" can be any number between 0 and 255. Such a block, where the first three numbers in the address are fixed and the last can be anything, is called a *class C subnet*. There are also class B subnets, which have the form 141.211.xxx.yyy, and class A subnets, which have the form 141.xxx.yyy.zzz.

Since all the addresses in a class C subnet are usually allocated to the same organization, a reasonable way of coarse-graining the Internet's network structure is to group nodes into class C subnets. In most cases this will group

⁴This description applies to addresses as they appear in IP version 4, which is the most widely used version of the protocol. A new version, version 6, which uses longer addresses, is slowly gaining acceptance, but it has a long way to go before it becomes as popular as its predecessor. (IP versions 1, 2, 3, and 5 were all experimental and were never used widely. Versions 4 and 6 are the only two that have seen widespread use.)

together nodes in the same organization, although larger organizations, like the University of Michigan, may own more than one class C subnet, so there will still be more than one node in the coarse-grained network corresponding to such organizations.

Given the topology of the network in terms of individual IP addresses, it is an easy matter to lump together into a single node all addresses in each class C subnet and place an edge between any two subnets if any address in one has a network connection to any address in the other. Figure 1.1 on page 2 shows an example of the network structure of the Internet at the level of class C subnets.

The second common type of coarse-graining is coarse-graining at the domain level. A *domain* is a group of computers and routers under, usually, the control of a single organization and identified by a single *domain name*, normally the last two or three parts of a computer's address when the address is written in human-readable text form (as opposed to the numeric IP addresses considered above). For example, "umich.edu" is the domain name for the University of Michigan and "oup.com" is the domain name for Oxford University Press. The name of the domain to which a computer belongs can be determined from the computer's IP address by a "reverse DNS lookup," a network service set up to provide precisely this type of information. Thus, given the network topology in terms of IP addresses, it is a straightforward task to determine the domain to which each IP address belongs and group nodes in the network according to their domain. Then an edge is placed between two nodes if any IP address in one has a direct network connection to any address in the other. The study by Faloutsos *et al.* [168] mentioned earlier looked at this type of domain-level structure of the Internet as well as the router-level structure.

The third common coarse-graining of the network is coarse-graining at the level of autonomous systems. This type of coarse-graining, however, is not usually used with traceroute data but with data obtained using an alternative method based on BGP routing tables, for which it forms the most natural unit of representation. The BGP method and autonomous systems are discussed in the next section.

2.1.2 MEASURING INTERNET STRUCTURE USING ROUTING TABLES

Internet routers maintain *routing tables* that allow them to decide in which direction incoming packets should be sent to best reach their destination. Routing tables are constructed from information shared between routers using BGP. They consist of lists of complete paths from the router in question to destinations on the Internet. When a packet arrives at a router, the router examines it to determine its destination and looks up that destination in the routing table.

The first step of the path in the appropriate table entry tells the router how the packet should be sent on its way.

In theory routers need store only the first step on each path in order to route packets correctly. However, for efficient calculation of routes using BGP it is highly desirable that routers be aware of the entire path to each destination, and since the earliest days of the Internet all routers have operated in this way. We can make use of this fact to measure the structure of the Internet.

Routing tables in routers are represented at the level of *autonomous systems*. An autonomous system (or AS) is a collection of routers, computers, or other devices, usually under single administrative control, within which data routing is handled independently of the wider Internet (hence the name “autonomous system”). That is, when a data packet arrives at a router belonging to an autonomous system, destined for a specific device or user within that same autonomous system, it is the responsibility of the autonomous system to get the packet the last few steps to its final destination. Data passing between autonomous systems, however, is handled by the Internet-wide mechanisms of BGP. Thus it’s necessary for BGP to know about routing only down to the level of autonomous systems and hence BGP tables are most conveniently represented in autonomous system terms. In practice, autonomous systems, of which there are (at the time of writing) about fifty thousand on the Internet, often coincide with domains, or nearly so.

Autonomous systems are assigned unique identification numbers. A routing path consists of a sequence of these AS numbers and since router tables contain paths to a large number of destinations we can construct a picture of the Internet at the autonomous system level by examining them. The process is similar to that for the traceroute method described in the previous section and depicted in Fig. 2.2. We must first obtain a set of router tables, which is normally done simply by asking router operators for access to their tables. Each router table contains a large number of paths starting from a single source (the router), and the union of the paths from many routers gives a good, though not complete, network snapshot in which the nodes are autonomous systems and the edges are the connections between autonomous systems. As with trace-route, it is important that the routers used be widely distributed across the network to avoid too much duplication of results, and the number of routers should be as large as possible to make the sampling of network edges as complete as possible. For example, the Routeviews Project,⁵ a large BGP-based Internet mapping effort based at the University of Oregon, uses (again at the

⁵See <http://www.routeviews.org>

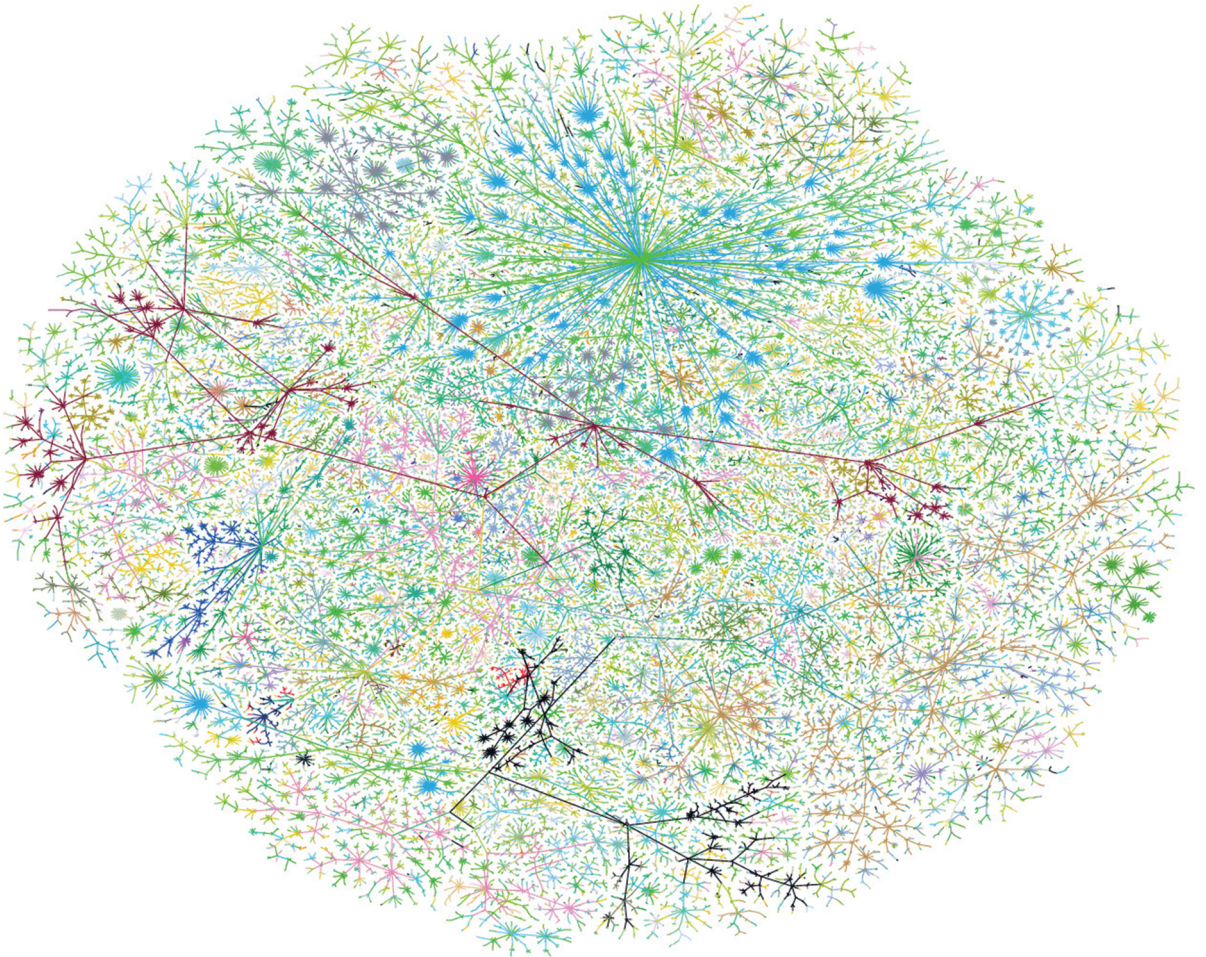


Figure 2.3: The structure of the Internet at the level of autonomous systems. The nodes in this network representation of the Internet are autonomous systems and the edges show the routes taken by data traveling between them. This figure is different from Fig. 1.1, which shows the network at the level of class C subnets. The picture was created by Hal Burch and Bill Cheswick. Patent(s) pending and Copyright Lumeta Corporation 2009. Reproduced with permission.

time of writing) a total of 501 source computers in 340 ASes around the world to measure the structure of the entire network every two hours.

Figure 2.3 shows a picture of the Internet at the AS level derived from routing tables. Qualitatively, the picture is similar to Fig. 1.1 for the class C subnet structure, but there are differences arising because class C subnets are smaller units than many autonomous systems and so Fig. 1.1 is effectively a

finer-grained representation than Fig. 2.3.

Using router-, subnet-, domain-, or AS-level structural data for the Internet, many intriguing features of the network's topology have been discovered in recent years [85, 102, 168, 323, 381, 384], some of which are discussed in later chapters of this book.

One further aspect of the Internet worth mentioning here is the geographic location of its nodes on the surface of the Earth. In many of the networks that we will study in this book, nodes do not exist at any particular position in real space—the nodes of a citation network, for instance, are not located on any particular continent or in any particular town. The nodes of the Internet, however, are by and large quite well localized in space. Your computer sits on your desk, a router sits in the basement of an office building, and so forth. Some nodes do move around, such as those representing mobile phones, but even these have a well-defined geographic location at any given moment. Things become a bit more blurry once the network is coarse-grained. The domain `umich.edu` covers large parts of the state of Michigan. The domain `ao1.com` covers most of North America. These are somewhat special cases, however, being unusually large domains. The majority of domains have a well-defined location at least to within a few miles. Furthermore, tools now exist for determining, at least approximately, the geographic location of a given IP address, domain, or autonomous system. Examples include *NetAcuity*, *IP2Location*, *MaxMind*, and many others. Geographic locations are determined primarily by looking them up in one of several registries that record the official addresses of the registered owners of IP addresses, domains, or autonomous systems. These addresses need not in all cases match the actual location of the corresponding computer hardware. For instance, the domain `ibm.com` is registered in New York City, but IBM's principal operations are in California. Nonetheless, an approximate picture of the geographic distribution of the Internet can be derived by these methods, and there has been some interest in the results [477].

For a review of work on geographic networks of various kinds see Barthélemy [46].

Geographic placement of nodes is a feature the Internet shares with several other technological networks, as we will see in the following sections, but rarely with networks of other kinds.⁶

⁶Social networks are perhaps the main exception. In many cases people or groups have reasonably well-defined geographic locations and a number of studies have looked at how geography and network structure interact [285, 300, 374, 439].

2.2 THE TELEPHONE NETWORK

The Internet is the best studied example of a technological network, at least as measured by the volume of recent academic work. This is partly because data on Internet structure are relatively easy to come by and partly because of intense interest among engineers, computer scientists, and the public at large. Other technological networks, however, are also of interest, including the telephone network and various distribution and transportation networks, and we look at some of these in the remainder of this chapter. Networks such as software call graphs and electronic circuits could also be considered technological networks and have been studied occasionally [174, 199, 334, 343, 485], but are beyond the scope of this book.

The telephone network—meaning the network of landlines and wireless links⁷ that transmits telephone calls—is one of the oldest electronic communication networks still in use, but it has been studied relatively little by network scientists, primarily because of a lack of good data about its structure. The structure of the phone network is known in principle, but the data are largely proprietary to the telephone companies that operate the network and, while not precisely secret, they are not openly shared with the research community in the same way that Internet data are. We hope that this situation will change, although the issue may become moot in the not too distant future, as telephone companies are sending an increasing amount of voice traffic over the Internet rather than over dedicated telephone lines, and it may not be long before the two networks merge into one.

Some general principles of operation of the telephone network are clear however. By contrast with the Internet, the traditional telephone network is not a packet-switched network of the kind described in Section 2.1. Signals sent over the phone network are not disassembled and sent as sets of discrete packets the way Internet data are (though there are exceptions—see below). The telephone network is a *circuit-switched* network, which means that the telephone company has a number of lines or circuits available to carry telephone calls between different points and it assigns them to individual callers when those

⁷For most of its existence, the telephone network has connected together stationary telephones in fixed locations such as houses and offices using landlines. Starting in the 1980s, fixed telephones have been replaced by wireless phones (“mobile phones” or “cell phones”), but it is important to realize that even calls made on wireless phones are still primarily carried over traditional landline networks. The signal from a wireless phone makes the first step of its journey wirelessly to a nearby transmission tower, but from there it travels over ordinary phone lines. Thus, while the advent of wireless phones has had an extraordinary impact on society, it has had rather less impact on the nature of the telephone network.

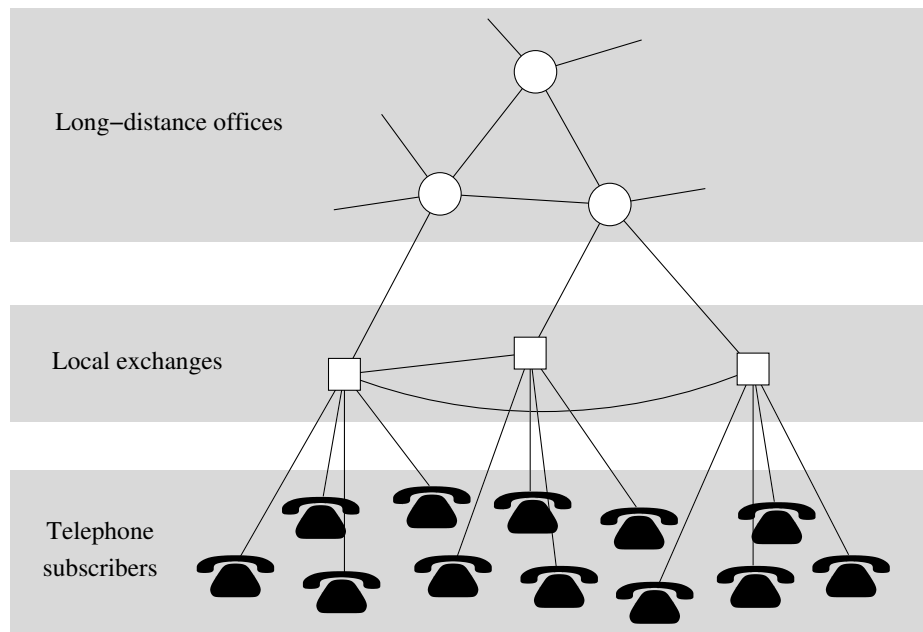


Figure 2.4: A sketch of the three-tiered structure of a traditional telephone network. Individual subscriber telephones are connected to local exchanges, which are connected in turn to long-distance offices. The long-distance offices are connected among themselves by trunk lines, and there may be some connections between local exchanges as well.

callers place phone calls. In the earliest days of telephone systems in the United States and Europe the “lines” actually were individual wires, one for each call the company could carry. Increasing the capacity of the network to carry more calls meant putting in more wires. Since the early part of the twentieth century, however, phone companies have employed techniques for *multiplexing* phone signals, i.e., sending many calls down the same wire simultaneously. The exception is the “last mile” of connection to the individual subscriber. The phone cable entering a house usually only carries one phone call at a time, although even that has changed in recent years as new technology has made it possible for households to have more than one telephone number and place more than one call at a time.

The basic form of the telephone network is relatively simple. Most countries with a mature landline telephone network use a three-tiered design, as shown in Fig. 2.4. Individual telephone subscribers are connected over local lines to

local telephone exchanges, which are then connected over shared “trunk” lines to long-distance offices, sometimes also called toll-switching offices. The long-distance offices are then connected among themselves by further trunk lines. The structure is, in many ways, rather similar to that of the Internet (Fig. 2.1), even though the underlying principles on which the two networks operate are different.

The three-level structure of the telephone network is designed to exploit the fact that most phone calls in most countries are local, meaning they connect subscribers in the same town or region. Phone calls between subscribers connected to the same local exchange can be handled by that exchange alone and do not need to make use of any trunk lines at all. Such calls are usually referred to as local calls, while calls that pass over trunk lines are referred to as trunk or long-distance calls. In many cases there may also be direct connections between nearby local exchanges that allow calls to be handled locally even when two subscribers are not technically attached to the same exchange.

The telephone network has had roughly this same topology for most of the past hundred years and still has it today, but many of the details about how the network works have changed. In particular, at the trunk level a lot of telephone networks are no longer circuit switched. Instead they are now digital packet-switched networks that work in a manner not dissimilar to the Internet, with voice calls being digitized, broken into packets, and transmitted over optical fiber links. Indeed, as mentioned, many calls are now transmitted digitally over the Internet itself, allowing phone companies to use the already existing Internet infrastructure rather than building their own. In many cases, only the “last mile” to the subscriber’s telephone is still carried on an old-fashioned dedicated circuit, and even that is changing with the advent of digital and Internet telephone services and mobile phones. Nonetheless, in terms of geometry and topology the structure of the phone network is much the same as it has always been, being dictated in large part by the constraints of geography and the propensity for people to talk more often to others in their geographic vicinity than to those further away.

2.3 POWER GRIDS

The power grid is the network of high-voltage transmission lines that provide long-distance transport of electric power within and between countries. The nodes in a power grid correspond to generating stations and switching substations, and the edges correspond to the high-voltage lines. (Low-voltage local power delivery lines are normally not considered part of the grid, at least where network studies are concerned.) The topology of power grids is not difficult to

determine. The networks are usually overseen by a single authority and complete maps of grids are readily available. Very comprehensive data on power grids (as well as other energy-related networks such as oil and gas pipelines) are available from specialist publishers, either on paper or in electronic form, if one is willing to pay for them.

There is much of interest to be learned by looking at the structure of power grids [13, 20, 31, 125, 263, 378, 415, 466]. Like the Internet, power grids have a spatial element; the individual nodes each have a location somewhere on the globe, and their distribution in space is interesting from geographic, social, and economic points of view. Network statistics, both geographic and topological, may provide insight into the global constraints governing the shape and growth of grids. Power grids also display some unusual behaviors, such as cascading failures, which can give rise to surprising outcomes such as the observed power-law distribution in the sizes of power outages [140, 263].

However, while there is a temptation to apply network models of the kind described in this book to try to explain the behavior of power grids, it is wise to be cautious. Power grids are complicated systems. The flow of power is governed not only by geometry and simple physical laws, but also by detailed control of the phases and voltages across transmission lines, monitored and adjusted on rapid timescales by sophisticated computer systems and on slower timescales by human operators. There is evidence to suggest that network topology has only a relatively weak effect on power failures and other power-grid phenomena, and that good prediction and modeling of power systems requires more detailed information than can be gleaned from a network representation alone [234, 378].

2.4 TRANSPORTATION NETWORKS

Another important class of technological networks are the transportation networks, such as airline routes and road and rail networks. The structure of these networks is not usually hard to determine. Airline networks can be reconstructed from published airline timetables, road and rail networks from maps. Geographic information systems (GIS) software can be useful for analyzing the geographic aspects of the data and there are also a variety of online resources providing useful information such as locations of airports.

One of the earliest examples of a study of a transportation network is the 1965 study by Pitts [387] of waterborne transport on Russian rivers in the Middle Ages. There was also a movement among geographers in the 1960s and 1970s to study road and rail networks, particularly focusing on the interplay between their physical structure and economics. The most prominent name

in the movement was that of Karel Karsky, and his book on transportation networks is a good point of entry into that body of literature [254].

More recently, a number of authors have produced studies applying new network analysis ideas to road, rail, air, and sea transportation networks [20, 34, 95, 198, 202, 224, 243, 290, 293, 324, 425, 426, 474]. In most of these studies the network nodes represent geographic locations and the edges represent routes. For instance, in studies of road networks the nodes usually represent road intersections and the edges roads. The study by Sen *et al.* [425] of the rail network of India provides an interesting counterexample. Sen *et al.* argue, plausibly, that in the context of rail travel what matters to most people is whether there is a direct train to their destination or, if there is not, how many trains they will have to take to get there. People do not care so much about how many stops there are along the way, so long as they don't have to change trains. Thus, Sen *et al.* argue, a useful network representation in the case of rail travel is one in which the nodes represent locations and two nodes are connected by an edge if a single train runs between them. Then the distance between two nodes in the network—the number of edges you need to traverse to get from A to B—is equal to the number of trains you would have to take. A better representation still (although Sen *et al.* did not consider it) would be a “bipartite network,” a network containing two types of node, one representing the locations and the other representing train routes. Edges in the network would then join locations to the routes that run through them. The first, simpler representation of Sen *et al.* can be derived from the bipartite one by making a “projection” onto the locations only. Bipartite networks and their projections are discussed in Section 6.6.

2.5 DELIVERY AND DISTRIBUTION NETWORKS

Falling somewhere between transportation networks and power grids are distribution networks, about which relatively little has been written within the field of networks research to date. Distribution networks include things like oil and gas pipelines, water and sewerage lines, and the routes used by the post office and package delivery companies. Figure 2.5 shows one example, the network of European gas pipelines, taken from a study by Carvalho *et al.* [96], who constructed the figure from data purchased from industry sources. In this network the edges are gas pipelines and the nodes are their intersections, including pumping, switching, and storage facilities and refineries.

If one is willing to interpret “distribution” in a loose sense, then one class of distribution networks that has been relatively well studied is river networks, though to be precise river networks are really collection networks rather than

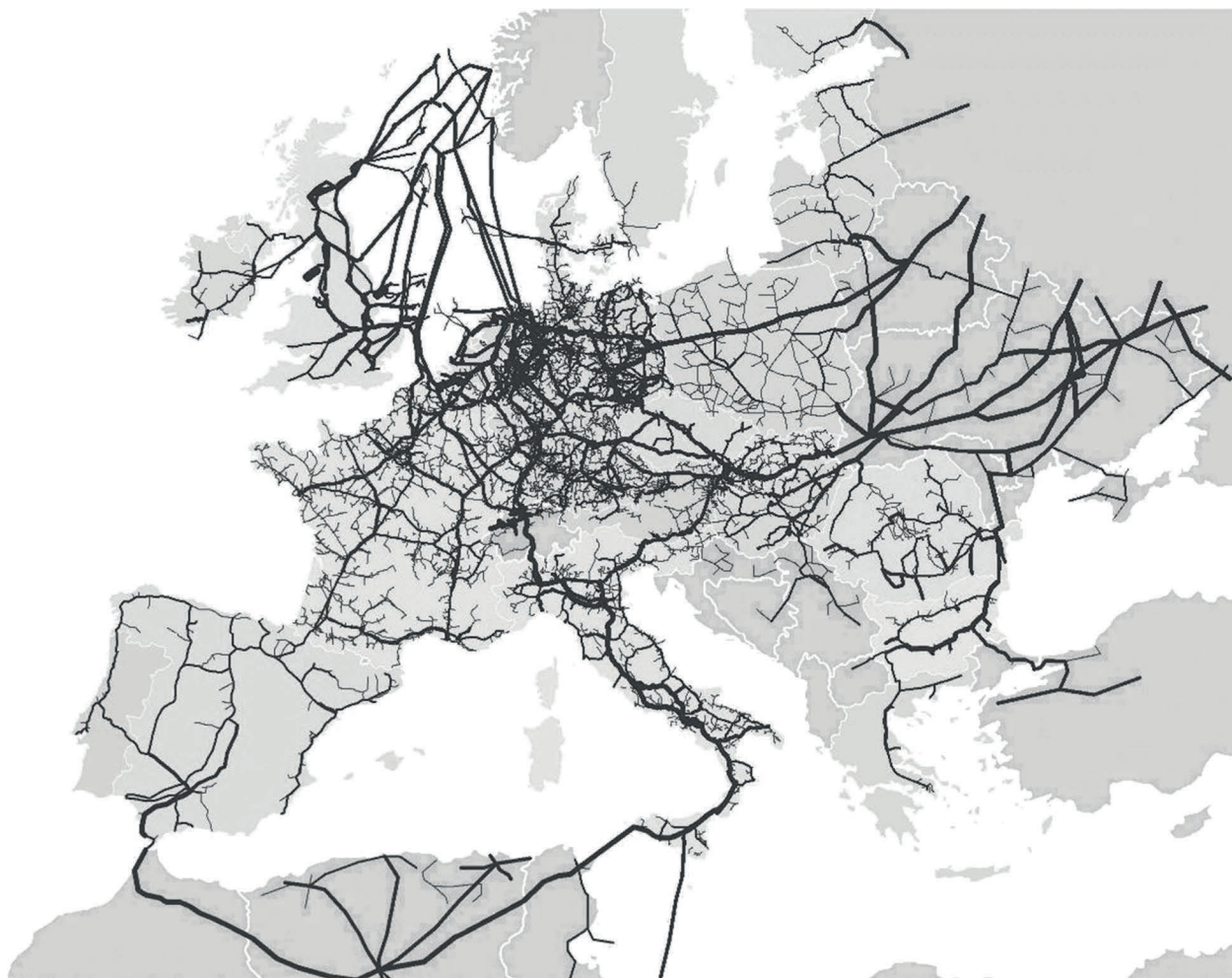


Figure 2.5: The network of natural gas pipelines in Europe. Thickness of lines indicates the sizes of the pipes. Reprinted with permission from R. Carvalho, L. Buzna, F. Bono, E. Gutierrez, W. Just, and D. Arrowsmith, Robustness of trans-European gas networks, *Phys. Rev. E* **80**, 016106 (2009). Copyright 2009 by the American Physical Society.

distribution networks. In a river network the edges are rivers or streams and the nodes are their intersections. As with road networks, no special techniques are necessary to gather data on the structure of river networks—the hard work of surveying the land has already been done for us by cartographers, and all we need do is copy the results from their maps. See Fig. 2.6 for an example.

The topological and geographic properties of river networks have been

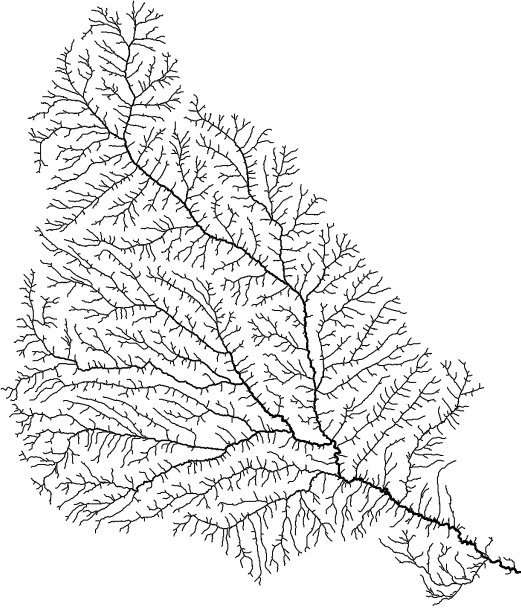


Figure 2.6: Drainage basin of the Loess Plateau. The network of rivers and streams on the Loess Plateau in the Shanxi province of China. The tree-like structure of the network is clearly visible—there are no loops in the network, so water at any point in the network drains off the plateau via a single path. Reproduced from Pelletier [386] by permission of the American Geophysical Union.

studied in some detail [143, 319, 407, 412]. Of particular note is the fact that river networks, to an excellent approximation, take the form of trees. That is, they contain no loops (if one disregards the occasional island midstream), a point that we discuss further in Section 6.8.

Similar in some respects to river networks are networks of blood vessels in animals, and their equivalents in plants, such as root networks. These too have been studied at some length. An early example of a mathematical result in this area is the formula for estimating the total geometric length of all edges in such a network by observing the number of times they intersect a regular array of straight lines [345]. This formula, whose derivation is related to the well-known “Buffon’s needle” experiment for determining the value of π , is most often applied to root systems, but there is no reason it could not also be useful in the study of river networks or, with suitable modification, any other type of geographic network.

Also of note in this area is work on the scaling relationships between the structure of branching vascular networks in organisms and metabolic processes [39, 468, 469], an impressive example of the way in which an understanding of network structure can be parlayed into an understanding of the functioning of the systems the networks represent. We will see many more examples during the course of this book.

CHAPTER 3

NETWORKS OF INFORMATION

A discussion of information networks, with a particular focus on the World Wide Web and citation networks

THIS CHAPTER focuses on networks of information, networks consisting of items of data linked together in some way. Information networks are all, so far as we know, man-made, with perhaps the best known example being the World Wide Web, though many others exist and are worthy of study, particularly citation networks of various kinds.

In addition, there are some networks which could be considered information networks but which also have social-network aspects. Examples include networks of email communications, networks on social-networking websites such as Facebook or LinkedIn, and networks of weblogs and online journals. We delay discussion of these and similar examples to the following chapter on social networks, in Section 4.4, but they could easily have fitted in the present chapter also. The classification of networks as information networks, social networks, and so forth is a fuzzy one, and there are plenty of examples that, like these, straddle the boundaries.

3.1 THE WORLD WIDE WEB

Although by no means the first information network created, the World Wide Web is probably the example best known to most people and a good place to start our discussion in this chapter.

Networks, 2nd edition. Mark Newman, Oxford University Press (2018). © Mark Newman.
DOI: 10.1093/oso/9780198805090.001.0001

As described in Chapter 1, the Web is a network in which the nodes are web pages, containing text, pictures, or other information, and the edges are the hyperlinks that allow us to navigate from page to page. The Web should not be confused with the Internet (Section 2.1), which is the physical network of data connections between computers; the Web is a network of links between pages of information.

Since hyperlinks run in one direction only, the Web is a directed network. We can picture the network with an arrow on each edge indicating which way it runs. Some pairs of web pages may be connected by hyperlinks running in both directions, which can be represented by two directed edges, one in each direction. Figure 3.1 shows a picture of a small portion of the web network, representing the connections between a set of web pages on a single website.

The World Wide Web was invented in the 1980s by scientists at the CERN high-energy physics laboratory in Geneva as a means of exchanging information among themselves and their co-workers, but it rapidly became clear that its potential was much greater [244]. At that time there were several similar information systems competing for dominance of the rapidly growing Internet, but the Web won the battle, largely because its inventors decided to give away for free the software technologies on which it was based—the Hypertext Markup Language (HTML) used to specify the appearance of pages and the Hypertext Transport Protocol (HTTP) used to transmit pages over the Internet. The Web’s extraordinary rise is now a familiar part of history and most of us use its facilities at least occasionally, and in many cases daily. A crude estimate of the number of pages on the Web puts that number at around 50 billion at the time of the writing¹ and it is, almost certainly, the largest network that has been studied quantitatively by network scientists to date.

The structure of the Web can be measured using a *crawler*, a computer

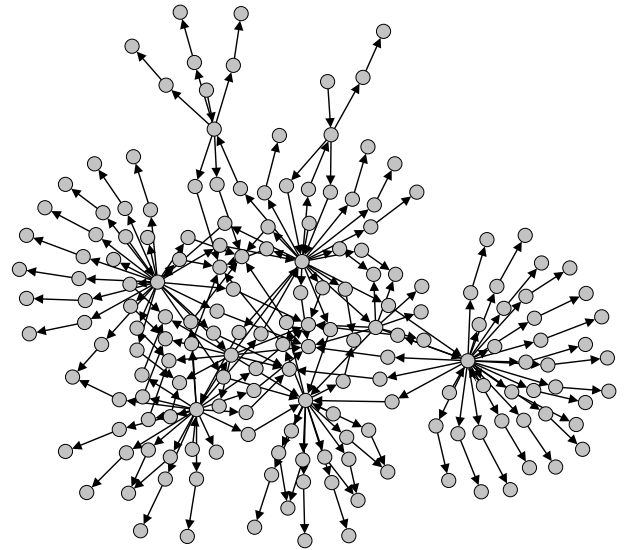


Figure 3.1: A network of pages on a corporate website. The nodes in this network represent pages on a website and the directed edges between them represent hyperlinks.

¹This is only the number of reachable static pages. The number of unreachable pages is difficult to estimate, and dynamic pages (see later) are essentially unlimited in number, although this may not be a very meaningful statement since these pages don’t exist until someone asks for them.

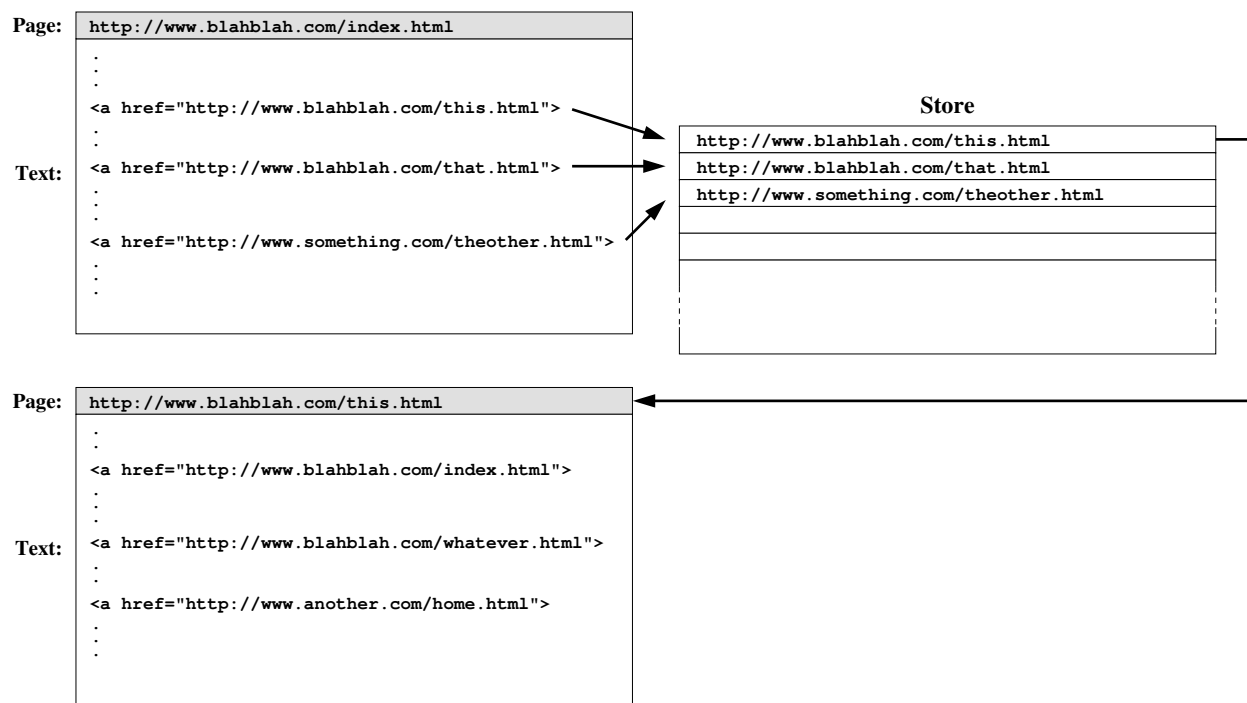


Figure 3.2: The operation of a web crawler. A web crawler iteratively downloads pages from the Web, starting from a given initial page. URLs are copied from the link tags in that initial page into a store. Once all links have been copied from the initial page, the crawler takes a URL from the store and downloads the corresponding page, then copies links from that, and so on.

Breadth-first search is discussed at length in Section 8.5.

program that automatically surfs the Web looking for pages. In its simplest form, the crawler performs a so-called breadth-first search on the web network, as shown schematically in Fig. 3.2. One starts from any initial web page, downloads the text of that page over the Internet, and finds all the links in the text. Functionally, a link consists of an identifying “tag”—a short piece of text marking the link as a link—and a *Uniform Resource Locator*, or URL, a standardized computer address that says how and where the linked web page can be found. By scanning for the tags and then copying the adjacent URLs a web crawler can rapidly extract URLs for all the links on a web page, storing them in memory or on a disk drive. When it is done with the current page, it takes one of the URLs from its store, uses it to locate a new page on the Web, and downloads the text of that page, and so the process repeats. If at any point the crawler encounters a URL it has seen before, then that URL is ignored and

not added to the store again, to avoid unnecessary duplication of effort. Only URLs that are different from those seen before are added to the store.

By repeating the process of downloading and URL extraction for a suitably long period of time, one can find a significant portion of the pages on the entire Web. No web crawler, however, finds all the pages on the Web, for a number of reasons. First, some websites forbid crawlers to examine their pages. Websites can place a file called `robots.txt` in their root directory that specifies which files, if any, crawlers can look at and may optionally specify that some crawlers are allowed to look at files while others are not. Compliance with the restrictions specified in a `robots.txt` file is voluntary, but in practice many crawlers do comply.

Second, many pages on the Web are dynamically generated: they are created on the fly by special software using, for instance, data from a database. Most large websites today, including many news, social media, retail, and corporate websites, as well as the web pages generated by search engines, fall into this category. Suppose, for instance, that you do a web search for “networks” using the Google search engine. Google does not keep a page of search results about networks (or anything else) just sitting on its computers, waiting for someone to ask for it. On the contrary, when you perform a search, the search engine rummages through its extensive database of web content (which it has found previously, using a web crawler) and makes a list of things that it believes will be useful to you. Then it creates a new web page containing that list and sends the page to your computer. The page of results you see when you search for something on Google is a dynamic page, generated automatically, and specifically for you, just a fraction of a second earlier.

As a result, the number of possible web pages that can be displayed as a result of a web search is so large as to be effectively infinite—as large as the number of different queries you could type into the search engine. When we are crawling the Web it is not practical for our crawler to visit all of these pages. The crawler must therefore make some choice about what it will look at and what it won’t. One choice would be to restrict ourselves to static web pages—ones that are not generated on the fly. But it’s not always simple to tell which pages are static, and besides, much useful information resides on the dynamic pages. In practice, the decisions made by crawlers about which pages to include tend to be fairly arbitrary, and it is not easy to guess which pages will be included in a crawl and which will not. But one can say with certainty that many will not and in this sense the crawl is always incomplete.

However, perhaps the most important reason why web crawls do not reach all the pages on the Web is that the network structure of the Web does not allow it. Since the Web is a directed network, not all pages are reachable from

a given starting point. In particular, it is clear that pages that have no incoming hyperlinks—pages that no one links to at all—can never be found by a crawler that follows links. Taking this idea one step further, it is also the case that a page will never be found if it is only linked to by pages that themselves have no incoming links. And so forth. In fact, the Web, and directed networks in general, have a special “component” structure, which we will examine in detail in Section 6.12.1, and most crawlers only find one part of that structure, the “giant out-component.” In the case of the World Wide Web the giant out-component is estimated to occupy only about a half of all web pages and the other half of the Web is unreachable [84].²

Although we are interested in web crawlers as a tool for probing the structure of the Web so that we can study its network properties, this is not their main purpose. The primary use of web crawlers is to construct directories of web pages for search purposes. Web search engines such as Google indulge in web crawling on a massive scale to find web pages and construct indexes of the words and pictures they contain that can later be used to locate pages of interest to searchers. Because their primary interest is indexing, rather than reconstructing the network structure of the Web, search engine companies don’t have any particular reason to take a good statistical sample of the Web and in network terms their crawls are probably quite biased. Still, many of them have graciously made their data available to academic researchers interested in web structure, and the data are good enough to give us a rough picture of what is going on. We will study a variety of features of the web network in subsequent chapters.

It isn’t entirely necessary that we rely on search engine companies or other web enterprises for data on the structure of the Web. One can also perform one’s own web crawls. There are a number of capable web crawler programs available for free on the Internet, including *wget*, *Nutch*, *GRUB*, and *Sphinx*. While most of us don’t have the time or the facilities to crawl billions of web pages, these programs can be useful for crawling small sets of pages or single websites, and much useful insight and information can be acquired by doing so.

²Which web pages a crawler finds does depend on where the crawl starts. A crawler *can* find a web page with no incoming links, for instance, if (and only if) it starts at that page. In practice, however, the starting point has remarkably little effect on what a crawler finds, since most of what is found consists of the giant out-component mentioned above, whose content does not depend on the starting point.

Web search, which itself raises some interesting network questions, is discussed in Section 18.1.

3.2 CITATION NETWORKS

A less well-known but much older information network is the network of citations between academic papers. Most papers reference one or more previous works, usually in a bibliography at the end of the paper, and one can construct a network in which the nodes are papers and there is a directed edge from paper A to paper B if A cites B in its bibliography. There are many reasons why one paper might cite another—to point out information that may be useful to the reader, to give credit for prior work, to highlight influences on the current work, or to disagree with the content of a paper. In general, however, if one paper cites another it is usually an indication that the contents of the earlier paper are relevant in some way to those of the later one, and hence citation networks are networks of relatedness of subject matter.

Quantitative studies of citation networks go back to the 1960s; perhaps the earliest is the 1965 study by Price [393]. Studies of citation networks fall within the field of information science, and more specifically within *bibliometrics*, the branch of information science that deals with the statistical study of publication patterns. The most common way to assemble a citation network is to do it by hand, simply typing the entries in the bibliographies of papers into a database from which a network can then be assembled. In the 1960s, when Price carried out his study, such databases were just starting to be created [200] and he made use of an early version of what would later become the Science Citation Index. Fifty years later, the Science Citation Index (along with its sister publications, the Social Science Citation Index and the Arts and Humanities Citation Index) is now one of the primary and most widely used sources of citation data. In recent years, it has moved from hand entry of bibliographic data to direct electronic submission of data by the journals, which makes for faster and more accurate database updates. Another database, Scopus, provides a competing but largely similar service. Both are professionally maintained and their coverage of the literature is reasonably complete and accurate, although the data are also quite expensive to purchase. Still, if one has the money, creating a citation network is only a matter of deciding which papers one wishes to include, using one of the databases to find the citations between those papers, and adding the appropriate directed edges to the network until it is complete.

More recently, software systems for compiling citation indexes automatically without human oversight have started to appear. Perhaps the best known of these is *Google Scholar*, the academic literature arm of the Google search engine. Google Scholar crawls the Web to find manuscripts of papers in electronic form and then searches through those manuscripts to identify citations to other papers. This is a somewhat hit-or-miss operation because many papers are not

Price's study is also the earliest we know of to find a power-law degree distribution in a network—see Section 10.4 for more discussion of this important phenomenon.

See Section 3.1 for a discussion of web crawlers.

on the Web or are not freely available, citations in papers have a wide variety of different formats and may include errors, and the same paper may exist in more than one place on the Web and possibly in more than one version. Nonetheless, enough progress has been made for Google Scholar to become a useful tool for the academic community. Other automated citation indexing projects include *Citebase*, which indexes physics papers, and *CiteseerX*, which indexes computer science.

As with web crawls, the original purpose of citation indexes was not to measure network structure. Citation indexes are assembled primarily to allow researchers to discover by whom a paper has been cited, and hence to find research related to a topic of interest. Nonetheless, data from citation indexes have been widely used to reconstruct the underlying networks and investigate their properties, and a number of large-scale studies of citation networks have appeared in recent years [101, 242, 294, 396–398, 404, 405].

Citation networks are in many ways similar to the World Wide Web. The nodes of the network hold information in the form of text and pictures, just as web pages do, and the links from one paper to another play a role similar to hyperlinks between web pages, alerting the reader when information relevant to the topic of one paper can be found in another. Papers with many citations are often more influential and widely read than those with few, just as is the case with web pages, and one can “surf” the citation network by following a succession of citations from paper to paper just as computer users surf the Web.

There is, however, at least one important difference between a citation network and the Web: a citation network is *acyclic*, while the Web is not. An acyclic network is one in which there are no closed loops of directed edges. On the World Wide Web, it is entirely possible to follow a succession of hyperlinks and end up back at the page you started at. On a citation network, by contrast, this is essentially impossible. The reason is that in order to cite a paper, that paper must already have been written. One cannot cite a paper that does not yet exist. Thus all the directed edges in a citation network point backward in time, from newer papers to older ones. If we follow a path of such edges from paper to paper, we will therefore find ourselves going backward in time, but there is no way to go forward again, so we cannot close the loop and return to where we started.³

Citation networks have some interesting statistics. For instance, one study

³On rare occasions it occurs that an author will publish two papers simultaneously in the same volume of a journal and, with the help of the printers, arrange for each paper to cite the other, creating a cycle of length two in the network. Thus the citation network is not strictly acyclic, having a small number of short loops scattered about it.

Academic studies of the Web within the information sciences sometimes refer to hyperlinks as “citations,” a nomenclature that emphasizes the close similarities between web and citation networks.

Acyclic networks are discussed further in Section 6.4.1.

See Fig. 6.3 for an illustration of a small acyclic network.

found that about 47% of all papers have never been cited at all [404]. Of the remainder, 9% have one citation, 6% have two, and it goes down quickly after that. Only 21% of all papers have 10 or more citations, and just 1% have 100 or more. These figures are a consequence of the power-law degree distribution of the network—see Section 10.4.

The most highly cited paper of all, according to the Science Citation Index, is a 1951 paper by Lowry *et al.* [311], which has been cited more than 300 000 times.⁴ Like most very highly cited papers, it is a methodological paper in molecular biology.

Citation networks of the type described so far are the simplest but not the only possible network representation of citation patterns. An alternative representation is the *cocitation network*. Two papers are said to be cocited if they are both cited by the same third paper. Cocitation is often taken as an indicator that papers deal with related topics and there is good evidence that this is a reasonable assumption in many cases. A cocitation network is a network in which the nodes represent papers and the edges represent cocitation of pairs of papers. By contrast with ordinary citation networks, the edges in a cocitation network are normally considered undirected, since cocitation is a symmetric relationship. One can also define a weighted cocitation network in which the edges have varying strengths: the strength of an edge between two papers is equal to the number of other papers that cite both.

Another related concept, although one that is less often used, is *bibliographic coupling*. Two papers are said to be bibliographically coupled if they cite the same other papers (rather than being cited by the same papers). Bibliographic coupling, like cocitation, can be taken as an indicator that papers deal with related material and one can define a strength or weight of coupling by the number of common citations between two papers. From the bibliographic coupling figures one can then assemble a bibliographic coupling network, either weighted or not, in which the nodes are papers and the undirected edges indicate bibliographic coupling.

3.2.1 PATENT AND LEGAL CITATIONS

The discussion of citation networks in the previous section focuses on citations between academic papers, but there are other types of citation also. Two of particular interest are citations between patents and between legal opinions.

Patents are temporary grants of ownership for inventions, which give their holders exclusive rights to control and profit from the protected inventions for a

⁴And it's been cited one more time now.

finite period of time. They are typically issued to inventors—either individuals or corporations—by national governments after a review process to determine whether the invention in question is original and has not been previously invented by someone else. In applying for a patent, an inventor must describe his or her invention in sufficient detail to make adequate review possible and present the case that the invention is worthy of patent protection. A part of this case typically involves detailing the relationship between the invention and other previously patented inventions, and in doing so the inventor will usually cite one or more previous patents. Citations may highlight dependencies between technologies, such as one invention relying for its operation on another, but more often patent citations are “defensive,” meaning that the inventor cites the patent for a related previous technology and then presents an argument for why the new technology is sufficiently different from the old one to merit its own patent. Governments, in the process of examining patent applications, will routinely consider their similarity to previous inventions, and defensive citations are one way in which an inventor can fend off in advance possible objections that might be raised. Typically, there are a number of rounds of communication back and forth between the government patent examiner and the inventor before a patent application is finally accepted or rejected. During this process extra citations are often added to the application, either by the inventor or by the examiner, to document the further points discussed in their communications.

If and when a patent is finally granted, it is published, citations and all, so that the public may know which technologies have patent protection. Published patents thus provide a source of citation data that we can use to construct networks similar to the networks of citations between papers. In patent networks the nodes are patents, each identified by a unique patent number, and the directed edges between them are citations of one patent by another. Like academic citation networks, patent networks are acyclic, or nearly so, with edges running from more recent patents to older ones, although short loops can arise in the network in the not uncommon case that an inventor simultaneously patents a number of mutually dependent technologies.

The structure of patent networks reflects the organization of human technology in much the same way that the structure of academic citation networks reflects the organization of research knowledge. Patent citations have been less thoroughly studied than academic citations, but the number of studies has been growing in the past few years with the appearance of high-quality data sources, including US National Bureau of Economic Research database

of US patents⁵ and the Google Patents search engine for worldwide patents.⁶ There are a number of interesting technological and legal questions, for instance concerning originality of patented inventions, emerging technologies, and antitrust policy, that can be addressed by examining patent citation networks [106, 161, 216, 247].

Another class of citation networks that have begun to attract attention in recent years are legal citation networks. In countries where law cases can be decided by judges rather than juries, such as civil cases or appeals in Europe or the US, a judge will frequently issue an “opinion” after deciding a case, a narrative essay explaining his or her reasoning and conclusions. It is common practice in writing such an opinion to cite previous opinions issued in other cases in order to establish precedent, or occasionally to argue against it. Thus, like academic papers and patents, legal opinions form a citation network, with opinions being the nodes and citations being the directed edges, and again the network is approximately acyclic. The legal profession has long maintained indexes of citations between opinions for use by lawyers, judges, scholars, and others, and in recent years these indexes have made the jump to electronic form and are now available online. In the United States, for instance, two commercial services, LexisNexis and Westlaw,⁷ provide detailed data on legal opinions and their citations. In the past few years a number of studies of the structure of legal citation networks have been published using data derived from these services [186, 187, 295, 314].

In principle it would be possible also to construct networks of cocitation or bibliographic coupling between either patents or legal opinions, but we are not aware of any studies yet published of such networks.

3.3 OTHER INFORMATION NETWORKS

There are many other kinds of information networks, although none have attracted the same level of attention as the Web and citation networks. In the remainder of this chapter we briefly discuss a few examples of other networks.

⁵See <http://www.nber.org/patents>

⁶See <http://patents.google.com>

⁷Westlaw is owned and operated by Thomson Reuters, the same company that owns the Science Citation Index, while LexisNexis is owned by Elsevier, which also owns Scopus.

3.3.1 PEER-TO-PEER NETWORKS

Peer-to-peer file-sharing networks (sometimes abbreviated P2P) are a widely used form of computer network that combines aspects of information networks and technological networks. A peer-to-peer network is a network in which the nodes are computers containing information in the form, usually, of discrete files, and the edges between them are virtual links established for the purpose of sharing the contents of those files. The links exist only in software—they indicate only the intention of one computer to communicate with another should the need arise.

Peer-to-peer networks are typically used as a vehicle for distributed databases, particularly for the storage and distribution, often illegally, of music and movies, although there are substantial legal uses as well, such as local sharing of files on corporate networks or the distribution of software. (The network of router-to-router communications using the Border Gateway Protocol described in Section 2.1 is another less obvious example of a legitimate and useful peer-to-peer network.)

The point of a peer-to-peer network is to facilitate the direct transfer of data between computers belonging to two end users of the network, two “peers.” This contrasts with the more common server–client model, such as that used by the World Wide Web, in which central server computers supply requested data to a large number of client machines. The peer-to-peer model is favored particularly for illicit sharing of copyrighted material because the owners of a centralized server can easily be obliged to deactivate the server by legal or law-enforcement action, but such actions are much more difficult when no central server exists. Eliminating central servers and the high-bandwidth connections they require also makes peer-to-peer networks economically attractive in applications such as software distribution.

On most peer-to-peer networks every computer is home to some information, but no computer has all the information in the network. If the user of a computer requires information stored on another computer, that information can be transmitted simply and directly over the Internet or over a local area network. This is a peer-to-peer transfer. No special infrastructure is necessary to accomplish it—standard Internet protocols are perfectly adequate to the task. Things get interesting, however, when one wants to *find* which other computer has the desired information. One way to do that is to have a central server containing an index of which information is on which computers. Such a system was employed by the early file-sharing network *Napster*, but a central server is, as we have said, susceptible to legal and other challenges, and such challenges

were in the end responsible for shutting Napster down.⁸

To avoid this problem, developers have turned to distributed schemes for searching and this is where network concepts come into play. In the simplest incarnation of the idea, computers form links to some number of their peers in such a way that all the computers together form a connected network. Again, a link here is purely a software construct—a computer’s network neighbors in the peer-to-peer sense are merely those others with which it has agreed to communicate when the need arises.

When a user instructs his or her computer to search the network for a specific file, the computer sends out a message to its network neighbors asking whether they have that file. If they do, they arrange to transmit it back to the user. If they do not, they pass the message on to *their* neighbors, and so forth until the file is found. As we show in Section 18.2, where we discuss search strategies on peer-to-peer networks at some length, this algorithm works, but only on relatively small networks. Since it requires messages to be passed between many computers for each individual search, the algorithm does not scale well as the network becomes large, the volume of network traffic generated by searches eventually swamping the available data bandwidth. To get around this problem, modern peer-to-peer networks employ a two-tiered network topology of nodes and “supernodes,” in which searches are performed only among the supernodes and ordinary nodes contact them directly to request searches be performed. More details are given in Section 18.2.

So what is the structure of a peer-to-peer network like? In many cases, unfortunately, not a lot is known since the software is proprietary and its owners are reluctant to share operational details. There have been a number of studies published of the early peer-to-peer network *Gnutella*, which was based on open-source software, meaning that the computer code for the software and the specification of the protocols it uses are freely available. By exploiting certain details of those protocols, particularly the ability for computers in the Gnutella network to “ping” one another (i.e., ask each other to identify themselves), researchers have been able to discover and analyze the structure of Gnutella networks [409, 442]. The networks appear to have approximately power-law degree distributions (see Section 10.4) and it has been suggested that this property could be exploited to improve search performance [6].

⁸The Napster name was later bought up by the music industry and is now the name of a legitimate online music service, although one that does not make use of peer-to-peer technology.

3.3.2 RECOMMENDER NETWORKS

Recommender networks represent people's preferences for things, such as for certain products sold by a retailer. Online merchants, for instance, may keep records of which customers bought which products and sometimes ask them whether they liked the products. Many large supermarket chains record the purchases made by their regular customers (usually identified by a small card with a barcode on it that is scanned when purchases are made) and so can work out which products each customer buys frequently.

We encountered bipartite networks previously in Section 2.4 and will study them further in Sections 4.5 and 6.6.

The fundamental representation of a recommender network is a “bipartite network,” a network with two types of node, one representing the products or other items and the other representing the people, with edges connecting people to the items they buy or like. One can also add strengths or weights to the edges to indicate, for instance, how often a person has bought an item or how much he or she likes it, or the strengths could be made negative to indicate dislikes.

Recommender networks have been studied for many types of goods and products, including books, music, films, and others. Interest in recommender networks arises primarily from their use in *collaborative filtering systems*, also sometimes called *recommender systems*, which are computer algorithms that attempt to guess new items a person will like by comparing their past preferences with those of other people. If person A likes many of the same things as person B, for instance, and if person B likes some further item that A has never expressed an opinion about, then maybe (the theory goes) A would like that item too. A wide variety of computer algorithms have been developed for extracting conclusions of this type from recommender networks [406] and are used extensively by retailers to suggest possible purchases to their customers, in the hope of drumming up business. The website of the online retailer *Amazon.com*, for instance, has a feature that recommends items to customers based on their previously expressed preferences and purchases. And many supermarkets now print out personalized discount coupons at checkout for products that a customer has not bought in the past but might be interested to try.

Product recommendations of this kind are big business: the ability to accurately predict what customers will like can mean millions of dollars in extra sales for a large retailer, or the difference between a loyal customer and one who defects to a competitor. In 2006, the entertainment company *Netflix* offered a prize of one million US dollars for anyone who could create a recommender system able to predict viewers' opinions about movies and TV programs 10% more accurately than the company's existing system. A mere 10% may not seem like a big improvement, but for a business the size of Netflix, with millions of

users, it could translate into a substantial increase in profits, easily justifying the prize money. Moreover, it turned out to be no trivial task to beat the 10% threshold. It took almost three years before the prize was finally won in 2009 by a large collaborative team of US and European researchers.

Research on recommender networks has focused mainly on the development of better collaborative filtering algorithms, but it is reasonable to suppose that the success of these algorithms should depend to some extent on the structure of the recommender network itself, and there is therefore good reason to also study that structure. A few such studies have been published in the scientific literature [94,220], but there is clearly room for further work.

3.3.3 KEYWORD INDEXES

Another type of information network, also bipartite in form, is the *keyword index*. An example is the index at the end of this book, which consists of a list of words or phrases, each accompanied by the numbers of the pages on which related information can be found. An index of this kind can be represented as a bipartite network, with two types of nodes representing words and pages, and an edge connecting each word to the pages on which it appears. In addition to their use in books, keyword indexes are routinely constructed as guides to other information collections, including sets of academic papers and the World Wide Web. The index constructed by a web search engine, as discussed in Section 3.1, is one example; it consists, at a minimum, of a set of words or phrases, with each word or phrase accompanied by a list of the web pages on which it occurs.

Indexes are of practical importance as a method for searching large bodies of information. Web search engines, for example, rely heavily on them to quickly find web pages that correspond to a particular query. However, indexes also have other, more sophisticated applications. They are used, for example, as a basis for techniques that attempt to find pages or documents that are similar to one another. Suppose one has a keyword index to a set of documents, consisting of a list of words and the documents they appear in. If we find that two documents contain a lot of the same keywords, it may be an indication that the two cover similar topics. A variety of computer algorithms for spotting such connections have been developed, typically making use of ideas very similar to those used in the recommender systems discussed in Section 3.3.2—the problem of finding documents with similar keywords is in many ways analogous to the problem of finding buyers who like similar products.

The identification of similar documents can be useful, for example, in constructing a search engine for searching through a body of knowledge. In a

standard index search, one typically looks up a keyword or set of keywords and gets a list of documents containing those words. Search engines that can tell when documents are similar may be able to respond more usefully because they can return documents that do not actually contain the keywords entered, but which are similar to documents that do. In cases where a single concept is called by more than one name, this may be an effective strategy for finding all the relevant documents.

In the context of document retrieval, the classic method for determining document similarity and performing generalized searches of this type is *latent semantic analysis*, which is based on the application of the matrix technique known as singular value decomposition to the bipartite network of keywords and documents [288]. A number of other competing methods have also been developed in recent years, using techniques such as non-negative matrix factorization [291, 292], latent Dirichlet allocation [63], and other probabilistic approaches [236].

As with recommender systems, it is reasonable to suppose that the success of methods for finding similar documents or improving searches using similarity information depends on the structure of the keyword index network, and hence that studies of that structure could generate useful insights. There has, however, been relatively little work on this problem so far within the network community, so there is plenty of room for future developments.

CHAPTER 4

SOCIAL NETWORKS

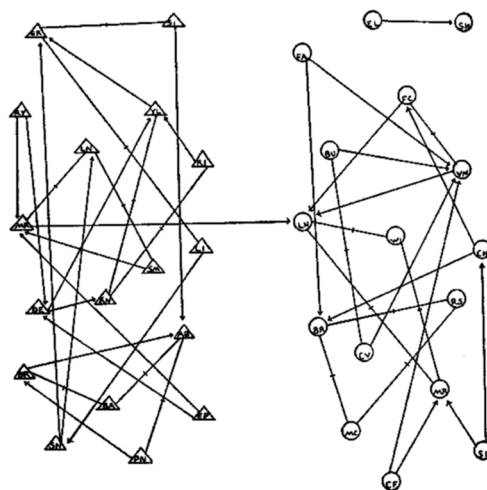
A discussion of social networks and the empirical techniques used to probe their structure

TO MOST people the words “social network,” if they mean anything, refer to online social media such as Facebook or Twitter. In the scientific study of networks, however, the phrase has a much broader meaning: a social network is any network in which the nodes represent people and the edges represent some form of connection between them, such as friendship. In this chapter we give a discussion of the origins and focus of the field of social network research and describe some of the types of networks studied and the techniques used to determine their structure. Sociologists have developed their own language for discussing social networks: they refer to the nodes, the people, as *actors* and the edges as *ties*. We will sometimes use these words when discussing social networks.

4.1 THE EMPIRICAL STUDY OF SOCIAL NETWORKS

Interest in social networks goes back many decades. Indeed, among researchers studying networks sociologists have perhaps the longest and best established tradition of quantitative, empirical work. There are clear antecedents of social network analysis to be found in the literature as far back as the end of the nineteenth century, though the true foundation of the field is usually attributed to psychiatrist Jacob Moreno, a Romanian immigrant to America who in the 1930s became interested in the dynamics of social interactions within groups of

Figure 4.1: Friendships between schoolchildren. This early hand-drawn image of a social network, taken from the work of psychiatrist Jacob Moreno, depicts friendship patterns between the boys (triangles) and girls (circles) in a class of schoolchildren in the 1930s. Reproduced from [341] by kind permission of the American Society of Group Psychotherapy and Psychodrama.



people. At a medical conference in New York City in March 1933 he presented the results of a set of investigations he had performed that may have been the first true social network studies, and the work attracted enough attention to merit a column in the *New York Times* a few days later. The following year Moreno published a book entitled *Who Shall Survive?* [341] which, though not a rigorous work by modern standards, contained the seeds of the field of *sociometry*, which later became social network analysis.

The most startling feature of Moreno's work was a set of hand-drawn figures depicting patterns of interaction among various groups of people. He called these figures *sociograms* rather than social networks (a term not coined until about twenty years later), but in everything but name they are clearly what we now know as networks. Figure 4.1, for instance, shows a diagram from Moreno's book, depicting friendships among a group of schoolchildren. The triangles and circles represent boys and girls respectively, and the figure reveals, among other things, that there are many friendships among the boys and many among the girls, but only one between a boy and a girl. It is simple conclusions like this, that are both sociologically interesting and easy to see once one draws a picture, that rapidly persuaded social scientists that there was merit in Moreno's methods.

One of the most important things to appreciate about social networks is that there are many different possible definitions of an edge in such a network and the particular definition one uses will depend on what questions one is interested in answering. Edges might represent friendship between individuals, but they could also represent professional relationships, exchange of goods

or money, communication patterns, romantic or sexual relationships, or many other types of connection. If one is interested, for instance, in professional interactions between the boards of directors of major corporations, then a network of who looks at who else's Facebook page is probably not of much use. Moreover, the techniques one uses to probe different types of social interaction can be quite different, so that very different kinds of studies may be needed to address different kinds of questions. Direct questioning of experimental subjects is probably the most common method of determining the structure of social networks. We discuss it in detail in Section 4.2.

Another important technique, the use of archival records, is illustrated by a different early example of a social network study. In 1939 a group of ethnographers studying the effects of social class and stratification in the American south collected data on the attendance of social events by 18 women in a small town in Mississippi over a period of nine months [129]. Rather than relying on interviews or surveys, however, they assembled their data using guest lists from the events and reports in the society pages of the newspapers.¹ Their study, often referred to as the "Southern Women Study," has been widely discussed and analyzed in the networks literature in the decades since its first publication. The data can be represented as a network in which the nodes represent the women and two women are connected if they attended a common event. An alternative and more complete representation is as an "affiliation network" or "bipartite network," in which there are two types of node representing the women and the events, and edges connecting each woman to the events she attended. A visualization of the affiliation network for the Southern Women Study is shown in Fig. 4.2.

One reason why this study has become so well known, in addition to its antiquity, is that the women were found by the researchers to split into two subgroups, tightly knit clusters of acquaintances with only rather loose between-cluster interaction. A classic problem in social network analysis is to devise a method or algorithm that can discover and extract such clustering from raw network data, and quite a number of researchers have made use of the Southern Women data as a test case for algorithm development.

Such is the power of social network analysis that its techniques have, since the time of Moreno and Davis *et al.*, been applied to an extraordinary variety of different communities, issues, and problems [79]: friendship and acquaintance patterns in local communities and in the population at large [54,55,261,333,447,452] and among university students [446,479] and schoolchildren [169,338,400];

¹They did also conduct some interviews, and made use of direct reports of attendance by observers. See Freeman [190] for a detailed discussion.

The use of archival and third-party records to reconstruct social networks is discussed in detail in Sections 4.4 and 4.5.

We encountered bipartite networks previously in Sections 2.4 and 3.3.2, and will study them in more detail in Sections 4.5 and 6.6.

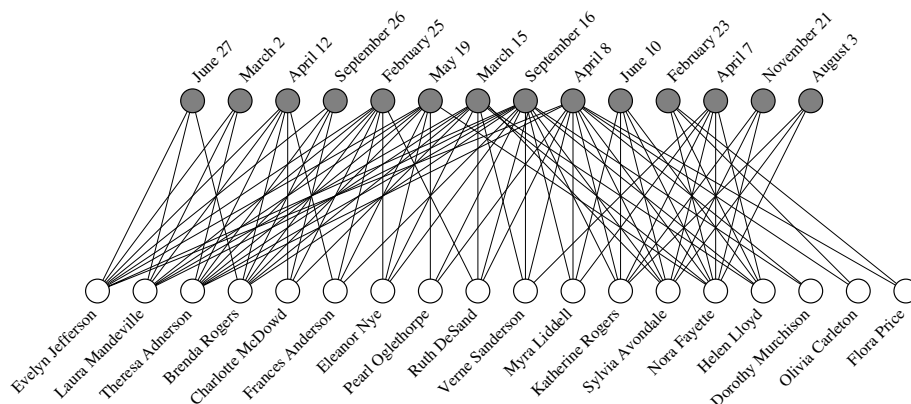


Figure 4.2: The affiliation network of the “Southern Women Study.” This network (like all affiliation networks) has two types of node, the open circles at the bottom representing the 18 women who were the subjects of the study and the shaded circles at the top representing the social events they attended. The edges connect each woman to the events she attended. Data courtesy of L. Freeman and originally from Davis *et al.* [129].

contacts between business people and other professionals [117, 197]; boards of directors of companies [130, 131, 318]; collaborations of scientists [218, 219, 349], movie actors [20, 466], and musicians [206]; sexual contact networks [271, 305, 392, 411, 417] and dating patterns [52, 238]; covert and criminal networks such as networks of drug users [421] or terrorists [282]; historical networks [51, 377]; online communities such as Usenet [313, 431, 449] and Facebook [278, 298, 446, 452]; and social networks of animals [180, 315, 418, 419].

We will see examples of these and other networks throughout this book and we will give details as needed as we go along. The rest of the present chapter is devoted to a discussion of the different empirical methods used to measure social networks. The techniques described above, direct questioning of subjects and the use of archival records, are two of the most important, but there are several others that find regular use. This chapter does not aim to give a complete review of the subject—for that we refer the reader to specialized texts such as those of Wasserman and Faust [462] and Scott [424]—but the material here provides a good grounding for our further studies in the remainder of the book.

4.2 INTERVIEWS AND QUESTIONNAIRES

The most common general method for gathering data on social networks is simply to ask people questions. If you are interested in friendship networks, you ask people who their friends are. If you are interested in business relationships you ask people who they do business with, and so forth. The asking may take the form of direct interviews with participants or the completion of questionnaires, either on paper or electronically. Indeed many modern studies, particularly telephone surveys, employ a combination of both interviews and questionnaires, wherein a professional interviewer reads questions from a questionnaire to a participant. By using a questionnaire, the designers of the study can guarantee that questions are asked, to a good approximation, in a consistent order and with consistent wording. By employing an interviewer to do the asking the study gains flexibility and reliability—interviewees often take studies more seriously when answering questions put to them by a human being—and interviewers may be given some latitude to probe interviewees when they are unclear, unresponsive, or confused. These are important considerations, since misunderstanding and inconsistent responses to survey questions are substantial sources of error [320]. By making questions as uniform as possible and giving respondents personal help in understanding them, these errors can be reduced. A good introduction to social survey design and implementation is given by Rea and Parker [403].

To measure social networks, surveys typically employ a *name generator*, a question or series of questions that invite respondents to name others with whom they have contact of a specific kind. For example, in their classic study of friendship networks among schoolchildren, Rapoport and Horvath [400] asked children to complete a questionnaire that included items worded as follows:

My best friend at ____ Junior High School is:
 My second-best friend at ____ Junior High School is:
 My third-best friend at ____ Junior High School is:
 ⋮
 My eighth-best friend at ____ Junior High School is:

The blanks “____” in the questionnaire were filled in with the appropriate school name.² The list stopped at the eighth-best friend and many participants did not complete all eight.

Ideally all students within a school would be surveyed, although Rapoport

²A junior high school in the United States is a school for children aged approximately 12 to 14 years.

and Horvath reported that in their case a few were absent on the day the survey was conducted. Note that the survey specifically asks children to name only friends within the school. The resulting network will therefore record friendship ties within the school but none to individuals elsewhere. This is a common issue: it is highly likely that any group of individuals surveyed will have at least some ties outside the group and one must decide what to do with these ties. Sometimes they are recorded. Sometimes, as here, they are not. Such details can be important since statistics derived from survey results will often depend on the decisions made.

There are a number of points to note about the data produced by name generators. First, the network ties, friendships in the case above, are determined by one respondent nominating another. This is a fundamentally asymmetric process. Individual A identifies individual B as their friend. In many cases B will also identify A as *their* friend, but there is no guarantee that this will happen and it is not uncommon for nomination to go in only one direction. We normally think of friendship as a two-way relationship, but surveys suggest that this not always the case. As a result, data derived from name generators are often best represented as directed networks, networks in which edges run in a particular direction from one node to another. If two individuals nominate each other then we have two directed edges, one pointing in either direction. Each node in the network then has two degrees, an out-degree—the number of friends identified by the corresponding individual—and an in-degree—the number of others who identified the individual as a friend.

This brings us to another point about name generators. It is common, as in the example above, for the experimenter to place a limit on the number of names a respondent can give. In the study of Rapoport and Horvath this limit was eight. Studies that impose such a limit are called *fixed choice* studies. The alternative is a *free choice* study, which imposes no limit.

Limits are often imposed purely for practical purposes, to reduce the work of the experimenter. However, they may also help respondents understand what is required of them. In surveys of schoolchildren, for instance, there are some children who, when asked to name their friends, will patiently name all the other children in the entire school, even if there are hundreds of them. Such responses are not particularly helpful—almost certainly the children in question are employing a different definition of friendship from that employed by most of their peers and by the investigators.

However, limiting the number of responses is for most purposes undesirable. In particular, it clearly limits the out-degree of the nodes in the network, imposing an artificial and possibly unrealistic cut-off. As discussed in Chapter 1, an interesting property of many networks is the existence of hubs, rare

We encountered directed networks previously in Chapter 1, in our discussion of the World Wide Web, and they are discussed in more detail in Section 6.4.

Recall that the degree of a node is the number of connections it has—see Section 6.10 for a detailed discussion.

nodes of unusually high degree, which, despite being few in number, can sometimes have a dominant effect on the behavior of the network. By employing a name generator that artificially cuts off the degree, any information about the existence of such hubs is lost.

It is worth noting, however, that even in a fixed choice study there is normally no limit on the *in*-degree of nodes in the network; there is no limit to the number of times an individual can be nominated by others. And indeed in many networks it is found that a small number of individuals are nominated an unusually large number of times. Rapoport and Horvath [400] observed this in their friendship networks: while most children in a school are nominated as a friend of only a few others, a small number of popular children are nominated very many times. Rapoport and Horvath were some of the first scientists in any field to study quantitatively the degree distributions of networks, reporting and commenting extensively on the *in*-degrees in their friendship networks.

Not all surveys employing name generators produce directed networks. Sometimes we are interested in ties that are intrinsically symmetric between the two parties involved, in which case the edges in the network are properly represented as undirected. An example is networks of sexual contact, which are widely studied to help us understand the spread of sexually transmitted diseases [271,305,392,417]. In such networks a tie between individuals A and B means that A and B had sex. While participants in studies sometimes do not remember who they had sex with or may be unwilling to talk about it, it is at least in principle a straightforward yes-or-no question whether two people had sex, and the answer should not depend on which of the two you ask. In such networks therefore, ties are normally represented as undirected.

Surveys can and often do ask respondents not just to name those with whom they have ties but to describe the nature of those ties as well. For instance, questions may ask respondents to name people they both like and dislike, or to name those with whom they have certain types of contact, such as socializing together, working together, or asking for advice. For example, in a study of the social network of a group of medical doctors, Coleman *et al.* [117] asked respondents the following questions:

Who among your colleagues do you turn to most often for advice?

With whom do you most often discuss your cases in the course of an ordinary week?

Who are the friends among your colleagues whom you see most often socially?

The names of a maximum of three doctors could be given in response to each question. A survey such as this, which asks about several types of interactions,

If individuals' responses differ too often, it is a sign that one's data are unreliable. Thus one may be able to estimate the level of measurement error in the data by comparing responses.

Multilayer networks are discussed further in Section 6.7.

The common tendency of people to associate with others who are similar to themselves in some way is called “homophily” or “assortative mixing,” and we discuss it in detail in Sections 7.7 and 10.7.

Experimental error in network measurements is discussed in detail in Chapter 9.

effectively generates data on several different networks at once—the network of advice, the discussion network, and so forth—but all built upon the same set of nodes. Networks such as this are sometimes called “multilayer” or “multiplex” networks.

Surveys may also pose questions aimed at measuring the strength of ties, asking, for instance, how often people interact or for how long, and they may ask individuals to give a basic description of themselves: their age, income, education, and so forth. Some of the most interesting results of social network studies concern the extent to which people’s choice of whom they associate with reflects their own background and that of their associates. For instance, you might choose to socialize primarily with others of a similar age to yourself, but turn for advice to those who are older than you.

The main disadvantages of network studies based on direct questioning of participants are that they are first laborious and second inaccurate. The administering of interviews or questionnaires and the collation of the responses is a demanding job that has been only somewhat eased by the use of computers and online survey tools. As a result, most studies have been limited to a few tens or at most hundreds of respondents—the 34-node social network of Fig. 1.2 is a typical example. It is a rare study that contains more than a thousand actors, and studies such as the National Longitudinal Study of Adolescent Health,³ which compiled responses from over 90 000 participants, are very unusual and extraordinarily costly. Only a substantial public interest such as, in that case, the control of disease, can justify the expense of performing them.

Data based on direct questioning may also be affected by biases of various kinds. Answers given by respondents are always to some extent subjective. If you ask people who their friends are, for instance, different people will interpret “friend” in different ways and thus give different kinds of answers, despite the best efforts of investigators to pose questions and record the answers in a uniform fashion. This problem is not unique to network studies. Virtually all social surveys suffer from such problems and a large body of expertise concerning techniques for dealing with them has been developed [320, 403]. Nonetheless, one should bear in mind when dealing with any social network derived from interviews or questionnaires the possibility of experimental bias in the data.

³See <http://www.cpc.unc.edu/projects/addhealth>

4.2.1 EGO-CENTERED NETWORKS

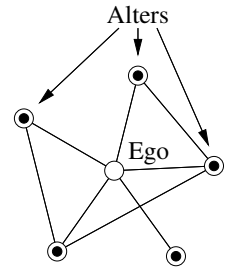
Studies in which all or nearly all of the individuals in a community are surveyed, as described in the previous section, are called *sociometric* studies, a term coined by Jacob Moreno himself (see the discussion at the beginning of this chapter). Sociometric studies are the gold standard for determining network structure but, as discussed at the end of the preceding section, they are also very labor intensive and for large populations may be infeasible.

At the other end of the spectrum from sociometric studies lie studies of *personal networks* or *ego-centered networks*.⁴ An ego-centered network is the network surrounding one particular individual, meaning that individual plus his or her immediate contacts. The individual in question is referred to as the *ego* and the contacts as *alters*.

Ego-centered networks are usually studied by direct questioning of participants, with interviews, questionnaires, or a combination of both being the instruments of choice (see Section 4.2). Typically, one constructs not just a single ego-centered network but several, centered on different egos drawn from the target population. In a telephone survey, for instance, one might call random telephone numbers in the target area and survey those who answer, asking them to identify others with whom they have a certain type of contact. Participants might also be asked to describe some characteristics both of themselves and of their alters, and perhaps to answer some other simple questions, such as which alters also have contact with one another.

Obviously, surveys of this type, and studies of ego-centered networks in general, cannot reveal the structure of an entire network. One receives snapshots of small local regions of the network, but in general those regions will not join together to form a complete social network. Sometimes, however, we are primarily interested in local network properties, and ego-centered network studies can give us good data about these. For example, if we wish to know about the degrees of nodes in a network—the numbers of ties people have—then a study in which a random sample of people are each asked to list their contacts may give us everything we need. (Studies probing node degrees are discussed more below.) If we also gather data on contacts between alters, we can estimate clustering coefficients (see Section 7.3). If we have data on characteristics of egos and alters we can measure assortative mixing (Sections 7.7 and 10.7).

An example of a study gathering ego-centered network data is the General



An ego-centered network consisting of an ego and five alters.

⁴Also called *egocentric* networks, although this term, which has its origins in social science and psychology, has taken on a different lay meaning which prompts us to avoid its use here.

Social Survey (GSS), a large-scale survey conducted every year in the United States starting in 1972 and every two years since 1994 [88]. The GSS is not primarily a social network study. The purpose of the study is to gather data about life in the United States, how it is changing, and how it differs from or relates to life in other societies. The GSS questionnaire contains a large number of parts, ranging from general questions probing the demographics and attitudes of the participants to specific questions about recent events, political topics, or quality of life. Among these many items, however, there are in each iteration of the survey a few questions about social networks. The precise number and wording of these questions changes from one year to another, but here are some examples from the survey of 1998, which was fairly typical:

From time to time, most people discuss important matters with other people. Looking back over the last six months, who are the people with whom you discussed matters important to you? Do you feel equally close to all these people?

Thinking now of close friends—not your husband or wife or partner or family members, but people you feel fairly close to—how many close friends would you say you have? How many of these close friends are people you work with now? How many of these close friends are your neighbors now?

By their nature these questions are of a “free choice” type, the number of friends or acquaintances the respondent can name being unlimited, although (and this is a criticism that has been leveled at the survey) they are also quite vague in their definition of friends and acquaintances, so people may give answers of widely varying kinds.

Another example of an ego-centered network study is the study by Bernard *et al.* [54, 55, 261, 326] of the degree of individuals in acquaintance networks (i.e., the number of people that people know). It is quite difficult to estimate how many people a person knows because most people cannot recall at will all those with whom they are acquainted and there is besides a lot of variation in people’s subjective definition of “knowing.” Bernard and co-workers came up with an elegant experimental technique for circumventing these difficulties. They asked study participants to read through a list of several hundred family names drawn at random from a telephone directory, and to count up how many people they knew with names appearing on the list. Each person with a listed name was counted separately, so that two acquaintances called Smith would count as two people. They were instructed to use the following precise definition of acquaintance:

Some care must be taken to match the selection of names to the community surveyed, since the frequency of occurrence of names shows considerable geographic and cultural variation.

You know the person and they know you by sight or by name; you can contact them in person by telephone or by mail; and you have had contact with the person in the past two years.

(Of course, many other definitions are possible. By varying the definition, one could probe different social networks.) Bernard and co-workers then multiplied the counts reported by participants by a scaling factor to estimate the total number of acquaintances of each participant. For instance, if the random names used in the study accounted for 1% of the population, then one would multiply by 100 to estimate the total number of acquaintances.

Bernard and co-workers repeated their study with populations drawn from several different US cities and the results varied somewhat from city to city, but overall they found that the typical number of acquaintances, in the sense defined above, of the average person in the United States is about 2000. In the city of Jacksonville, Florida, for instance, they found a figure of 1700, while in Orange County, California they found a figure of 2025. Many people find these numbers surprisingly high upon first encountering them, perhaps precisely because we are poor at recalling all of the many people we know. But repeated studies have confirmed figures of the same order of magnitude, at least in the United States. In some other countries the figures are different. Bernard and co-workers repeated their study in Mexico City, for instance, and found that the average person there knows about 570 others.

4.3 DIRECT OBSERVATION

An obvious method for constructing social networks is direct observation. Simply by watching interactions between individuals one can, over a period of time, form a picture of the networks of unseen ties that exist between those individuals. Most of us, for instance, will be at least somewhat aware of friendships or enmities that exist between our friends or co-workers. In direct observation studies, researchers attempt to develop similar insights about whatever population they are interested in.

Direct observation tends to be a labor-intensive method of study, so its use is usually restricted to small groups, and primarily ones with extensive face-to-face interactions in public settings. In Chapter 1 we saw one example, the “karate club” network of Zachary [479] (see Fig. 1.2 on page 5). Another is the study by Freeman *et al.* [193, 194] of the social interactions of a group of windsurfers, in which experimenters watched windsurfers on a beach in Orange County, California and recorded the length in minutes of every pairwise interaction among them. A large number of direct-observation network data

sets were compiled by Bernard and co-workers during the 1970s and 1980s as part of a lengthy study of the accuracy of individuals' perception of their own social situation [56, 58, 59, 259]. These included data sets on interactions among students, faculty, and staff in a university department, on members of a university fraternity,⁵ on users of a teletype service for the deaf, and several other examples.

One arena in which direct observation is essentially the only viable experimental technique is studies of the social networks of animals, since clearly animals cannot be surveyed using interviews or questionnaires. Not all animal species form interesting social networks, but informative studies have been performed of, among others, monkeys [180, 418, 419], kangaroos [214], and dolphins [121, 315]. A common approach is to record instances of animal pairs engaging in recognizable social behaviors such as mutual grooming, courting, or close association, and then to declare ties to exist between the pairs that engage in these behaviors most often. Networks in which the ties represent aggressive behaviors have also been reported, such as networks of baboons [328], bison [310], deer [27], wolves [249, 455], and ants [116]. In cases where aggressive behaviors normally result in one animal's establishing dominance over another the resulting networks can be regarded as directed and are sometimes called *dominance hierarchies* [136, 137, 150].

4.4 DATA FROM ARCHIVAL OR THIRD-PARTY RECORDS

An increasingly important, voluminous, and often highly reliable source of social network data is archival records. Such records are, at least sometimes, relatively free from the vagaries of human memory and can be impressive in their scale, allowing us to construct networks of a size that would be unreachable by other methods. Archival records can also allow us to reconstruct networks that no longer exist, such as networks from the historical past.

A well-known, small-scale example of a study based on archival records is the work of Padgett and Ansell on the ruling families of Florence in the fifteenth century [377]. In this study the investigators looked at contemporaneous historical records to determine which among the Florentine families had trade relations, marriage ties, or other forms of social contact with one another. Figure 4.3 shows one of the resulting networks, a network of intermarriages between 15 of the families. It is notable that the Medici family occupies a central

⁵In American universities a "fraternity" is a combined social organization and boarding house for male students.

position in this network, having marriage ties with members of no fewer than six other families, and Padgett and Ansell conjectured that it was by shrewd manipulation of social ties such as these that the Medici rose to a position of dominance in Florentine society.

In recent years, researchers have used archival records to construct a wide variety of different networks, some of them very large. A number of authors, for example, have looked at email networks [156, 277, 450]. Drawing on email logs—automatic records kept by email servers of messages sent and received—it is possible to construct networks in which the nodes are people (or more correctly email addresses) and the directed edges between them are email messages. Exchange of email in such a network can be taken as a proxy for acquaintance between individuals, or we may be interested in patterns of email exchange for some other reason, such as understanding how information spreads through a community. Similar networks can also be constructed from patterns of text messaging or instant messaging using mobile phones [374, 439].

A network similar in some ways the email network is the *telephone call graph*, in which the nodes represent telephone numbers and directed edges between them represent telephone calls from one number to another. Call graphs can be constructed from call logs kept by telephone companies, and a number of studies have been performed in recent years, including some at the largest scales, with a million or more phone numbers [1, 10, 64, 233, 375, 401]. Studies of mobile phones have attracted particular attention because mobile phone data can reveal not only who calls whom but also potentially the geographic location of the phone users, providing a rare opportunity to construct networks with both detailed contact patterns and high spatial resolution [285, 374, 439]. Mobile phone data have also played a role in studies of face-to-face social interaction: if two phones are recorded as being in the same location at the same time one can perhaps conclude that their owners had face-to-face contact, and a number of studies have been conducted using assumptions of this kind [91, 154, 155, 439].

Email networks and telephone call graphs have another feature of partic-

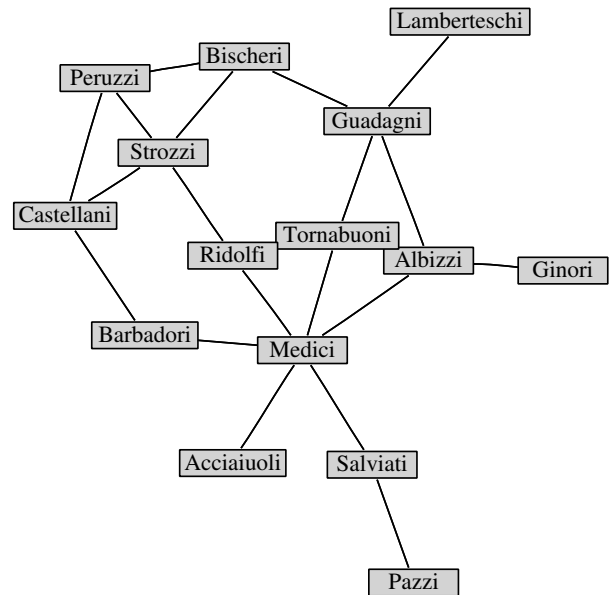


Figure 4.3: Intermarriage network of the ruling families of Florence in the fifteenth century. In this network the nodes represent families and the edges represent ties of marriage between them. After Padgett and Ansell [377].

Telephone call graphs are quite distinct from the physical network of telephone cables discussed in Section 2.2 and the two should not be confused. Indeed, a call graph is to the physical telephone network roughly as an email network is to the Internet.

In sociology, studies of time-varying networks are sometimes called *longitudinal* studies, and you may occasionally encounter this term in the literature.

ular interest: they are time-resolved—the date and time of each interaction is in principle known, allowing us to reconstruct after the fact the timing and duration of contacts between individuals if we have access to the appropriate data. Most of the sources of network data considered in this book are not time-resolved, but many networks do nonetheless change over time. Time-varying networks have been the focus of increasing research attention in recent years. We discuss them further in Section 6.7.

Recent years have also seen the rapid emergence of online social networking services, such as Facebook and LinkedIn, which, as a natural part of their operation, build records of connections between their participants and hence provide a rich source of archival network data. Some, such as Twitter, have made their data (or a part of it) publicly available, allowing researchers to study the corresponding networks [141, 208, 212]. Others are not publicly available but the companies involved have in some cases published analyses of their own networks, with some, such as Facebook, operating substantial internal research departments or inviting academic researchers to collaborate [28, 74, 278, 452]. Some online communities are not explicitly oriented towards networks or networking but can be studied using network techniques nonetheless. A number of researchers have looked, for instance, at networks of interactions between users of online dating sites [238, 297].

Weblogs, online diaries and journals, and other kinds of personal websites are another source of online social network data, although their popularity has waned somewhat in recent years. On these sites an individual or sometimes a group of people post their thoughts on topics of interest, often accompanied by links to other sites, and the sites and links form a directed network that lies, in terms of semantic content, somewhere between a social network and the World Wide Web: the links are often informational—the linker wishes to bring to his or her readers' attention the contents of the linked site—but there is a strong social element as well, since people often link to sites operated by their friends or acquaintances. The structure of the networks of links can be extracted using crawlers similar to those used to search the Web—see Section 3.1. Studies of weblogs and journals have been performed, for example, by Adamic and Glance [4] and MacKinnon and Warren [317].

4.5 AFFILIATION NETWORKS

An important special case of network data from archival records is the *affiliation network*. An affiliation network is a network in which actors are connected via their membership in groups of some kind. We saw one example at the start of this chapter, the Southern Women Study of Davis *et al.* [129], in which women

were connected via their common attendance at social events: the groups in that case were the attendees of the events. As we saw, the most complete representation of an affiliation network is a network with two types of nodes representing the actors and the groups, with edges connecting actors to the groups to which they belong—see Fig. 4.2 on page 50. In such a representation, called a “bipartite network” or “two-mode network,” there are no edges connecting actors directly to other actors or groups to other groups, only actors to groups.

Many examples of affiliation networks can be found in the literature. A famous case is the study by Galaskiewicz [197] of the CEOs of companies in Chicago in the 1970s and their social interaction via clubs that they attended: the CEOs are the actors and the clubs are the groups. Also in the business domain, a number of studies have been conducted of the networks formed by the boards of directors of companies [130,131,318], where the actors are company directors and the groups are the boards on which they sit. In addition to looking at the connections between directors in such networks, which arise as a result of their sitting on boards together, attention has also been focused on the connections between boards (and hence between companies) that arise as a result of their sharing a common director, a so-called board “interlock.”

More recently, some extremely large affiliation networks have been studied in the mathematics and physics literature. Perhaps the best known example is the network of collaboration of film actors, in which the “actors” in the network sense are actors in the dramatic sense also, and the groups to which they belong are the casts of films. This network is the basis, among other things, for a well-known parlor game, sometimes called the “Six Degrees of Kevin Bacon,” in which one attempts to connect pairs of actors via chains of intermediate costars, in a manner reminiscent of the small-world experiments of Stanley Milgram which we discuss in Section 4.6. The film actor network has, with the advent of the Internet, become very thoroughly documented and has attracted the attention of network analysts in recent years [20,40,466], although it is not clear whether there are any conclusions of real scientific interest to be drawn from its study.

Another example of a large affiliation network, one that holds more promise of providing useful results, is the coauthorship network of academics. In this network the actors are academic authors and the groups are the sets of authors of learned papers. Like the film actor network, this network is well documented, for instance via online bibliographic databases of published papers. Whether one is interested in papers published in journals or in more informal forums such as online preprint servers, excellent records now exist in most academic fields of authors and the papers they write, and a number of studies of the

We study bipartite networks in more detail in Section 6.6.

corresponding affiliation networks have been published [43, 133, 196, 218, 219, 241, 347–349, 460].

4.6 THE SMALL-WORLD EXPERIMENT

A memorable and illuminating contribution to the social networks literature was made by the psychologist Stanley Milgram in the 1960s with his now-famous “small-world” experiments [333, 447]. Milgram was interested in quantifying the typical distance between actors in social networks. As discussed in Chapter 1, one can define the distance between two nodes in a network as the number of edges that must be traversed to go from one node to the other. There are mathematical arguments that suggest that this distance should be small for most pairs of nodes in most networks (see Section 11.7), a fact that was already well known in Milgram’s time.⁶ Milgram wanted to test this conjecture under real-world conditions and to do this he concocted the following experiment.

Milgram conducted several sets of small-world experiments. The one described here is the first and most famous, but there were others—see Refs. [275, 447].

Milgram sent a set of packages, 96 in all, to volunteer participants in the US town of Omaha, Nebraska, who were recruited via a newspaper advertisement. The packages contained an official-looking booklet, or “passport,” emblazoned in gold letters with the name of Milgram’s home institution, Harvard University, plus a set of written instructions. The instructions asked the participants to get the passport to a specified target individual, a friend of Milgram’s who lived in Boston, Massachusetts, over a thousand miles away. The only information supplied about the target was his name (and hence indirectly the fact that he was male), his address, and his occupation as a stockbroker. But the passport holders were not allowed simply to send their passport to the given address. Instead they were asked to send it to someone they knew on a first-name basis, more specifically to the person in this category who they felt would stand the best chance of getting the passport to the intended target. Thus they might decide to send it to someone they knew who lived in Massachusetts, or maybe someone who worked in the financial industry. The choice was up to them. Whoever they did send the passport to was then to repeat the process, sending it to one of *their* acquaintances, and so forth, so that after a succession of steps the passport would, with luck, find its way into the hands of its intended recipient. Since every step of the process corresponded to the passport’s changing hands between a pair of first-name acquaintances, the entire journey consti-

⁶Milgram was particularly influenced in his work by a mathematical paper by Pool and Kochen [389] that dealt with the small-world phenomenon and had circulated in preprint form in the social science community for some years when Milgram started thinking about the problem, although the paper was not officially published until some years later.

tuted a path along the edges of the social network formed by the set of all such acquaintanceships, and the length of the journey provided an upper bound on the distance through this network between the starting and ending individuals in the chain.

Of the 96 passports sent out, 18 found their way to the stockbroker target in Boston. (While this may at first sound like a low figure, it is actually quite high—recent attempts to repeat Milgram’s work have resulted in response rates orders of magnitude lower [142].) Milgram asked participants to record in the passport each step of the path taken, so he knew how long each path was, and he found that the mean length of completed paths from Omaha to the target in Boston was just 5.9 steps. This result is the origin of the idea of the “six degrees of separation,” the popular belief that there are only about six steps between any two people in the world.

For a number of reasons this result is probably not very accurate. The initial recipients in the study were not chosen at random—they were volunteers who answered an advertisement—so they may not have been typical members of the population. At the very least, all of them were in a single town in a single country, which calls into question the extent to which the results of the study apply to the population of the world as a whole, or even to the population of the United States. Furthermore, Milgram used only a single target in Boston, and there is no guarantee this target was typical of the population either. Also we don’t know that chains took the shortest possible route to the target. Probably they did not, at least in some cases, so the lengths of the chains provide only an upper bound on the actual distance between nodes. Moreover, most of the chains were never completed. Many passports were discarded or lost and never made their way to the target. It is reasonable to suppose that the chances of getting lost were greater for passports that took longer paths, and hence that the paths that were completed were a biased sample, having typical lengths shorter than the average.

For all of these reasons Milgram’s results should be taken with a large pinch of salt. Even so, the fundamental conclusion that node pairs in social networks tend on average to be connected by short paths is now widely accepted. It has been confirmed directly in many cases, including for some very large social networks such as the entire network of Facebook friendships [452], and has moreover been shown to extend to many other (non-social) kinds of networks as well. Enough experiments have observed this “small-world effect” in enough networks that, whatever misgivings we may have about Milgram’s particular technique, the general result is not seriously called into question.

Milgram’s experiments also, as a bonus, revealed some other interesting features of acquaintance networks. For instance, Milgram found that most of

The phrase “six degrees of separation” did not appear in Milgram’s writing. It is more recent and comes from the title of a successful Broadway play by John Guare [221], later made into a film, in which the lead character discusses Milgram’s work.

the passports that did find their way to the stockbroker target did so via just three of the target's friends. That is, a large fraction of the target's connections to the outside world seemed to be through only a few of his acquaintances, a phenomenon sometimes referred to as "funneling." Milgram called such well-connected acquaintances "sociometric superstars," and their existence has occasionally been noted in other networks also, such as collaboration networks [347], although not in some others [142].

A further interesting corollary of Milgram's experiment, never mentioned by Milgram himself, was highlighted many years later by Kleinberg [266,267]: the fact that a moderate number of passports did find their way to the intended target person shows not only that short paths exist in the acquaintance network, but also that people are good at finding those paths. Upon reflection this is quite a surprising result. As Kleinberg has shown, it is possible and indeed common for a network to possess short paths between nodes but for them to be hard to find unless one has complete information about the structure of the entire network, which the participants in Milgram's studies did not. Kleinberg has conjectured that the network of acquaintances needs to have a special type of structure for the participants to find the paths they did with only limited knowledge of the network. We discuss his ideas in detail in Section 18.3.

Recently the small-world experiment has been repeated by Dodds *et al.* [142] using the modern medium of email. In this version of the experiment participants forwarded email messages to their acquaintances in an effort to get them to a specified target person about whom they were told a few basic facts. The experiment improved on Milgram's in terms of sheer volume, and also by having much more numerous and diverse target individuals and starting points for messages: 24 000 chains were started, most (though not all) with unique starting individuals, and with 18 different participating targets in 13 different countries. On the other hand, the experiment experienced enormously lower rates of participation than Milgram's, perhaps because the public is by now quite jaded in its attitude towards unsolicited mail. Of the 24 000 chains, only 384, or 1.5%, reached their intended targets, compared with 19% in Milgram's case. Still, the basic results were similar to those of Milgram. Completed chains had an average length of just over four steps. Because of their better data and more careful statistical analysis, Dodds *et al.* were also able to compensate for biases due to unfinished chains and estimated that the true average path length for the experiment was somewhere between five and seven steps—very similar to Milgram's result. However, Dodds *et al.* observed no equivalent of the "sociometric superstars" of Milgram's experiment, raising the question of whether their appearance in Milgram's case was a fluke of the particular target individual he chose rather than a generic property of social networks.

An interesting variant on the small-world experiment has been proposed by Killworth and Bernard [57,260], who were interested in how people “navigate” through social networks, and specifically how participants in the small-world experiments decided whom to forward messages to in the effort to reach the specified target. They conducted what they called “reverse small-world” experiments⁷ in which they asked participants to *imagine* that they were taking part in a small-world experiment. A (fictitious) message was to be communicated to a target individual and participants were asked what they would like to know about the target in order to decide whom to forward the message to. The actual passing of the message never took place; the experimenters merely recorded what questions participants asked about the target. They found that three characteristics were sought overwhelmingly more often than any others, namely the name of the target, their geographic location, and their occupation—the same three pieces of information that Milgram provided in his original experiment. Some other characteristics came up with moderate frequency, particularly when the experiment was conducted in non-Western cultures or among minorities: in some cultures, for instance, parentage or religion were considered important identifying characteristics of the target.

While the reverse small-world experiments do not directly tell us about the structure of social networks, they do give us information about how people perceive and deal with social networks.

4.7 SNOWBALL SAMPLING, CONTACT TRACING, AND RANDOM WALKS

Finally in this chapter on social networks we take a look at a class of network-based techniques for sampling hidden populations.

Studies of some populations, such as drug users or illegal immigrants, present special problems to the investigator because the members of these populations do not usually want to be found and are often wary of giving interviews. Techniques have been developed, however, for sampling these populations by making use of the social networks that connect their members together. The most widely used such technique is *snowball sampling* [162, 188, 445].

Note that, unlike the other experimental techniques discussed in this chapter, snowball sampling is not intended as a technique for probing the structure of social networks. Rather, it is a technique for studying hidden populations

⁷Also sometimes called INDEX experiments, which is an abbreviation for “informant-defined experiment.”

The mechanisms of network search and message passing are discussed in greater detail in Section 18.3.

that relies on social networks for its operation. It is important to keep this distinction clear. To judge by the literature, some professional network scientists do not, a mistake that can result in erroneous conclusions and bad science.

Standard techniques such as telephone surveys often do not work well when sampling hidden populations. An investigator calling a random telephone number and asking if anyone on the other end of the line uses drugs is unlikely to receive a useful answer. The target population in such cases is small, so the chances of finding one of its members by random search are slim, and when you do find one they will very likely be unwilling to discuss the highly personal and possibly illicit topic of the survey with an investigator they have never met before and have no reason to trust.

So investigators probe the population instead by getting some of its members to provide contact details for others. The typical survey starts off rather like a standard ego-centered network study (see Section 4.2.1). You find one initial member of the population of interest and interview them about themselves. Then, upon gaining their confidence, you invite them also to name other members of the target population with whom they are acquainted. Then you go and find those acquaintances and interview them in turn, asking them also to name further contacts, and so forth through a succession of “waves” of sampling. Pretty soon the process “snowballs” and you have a large sample of your target population to work with.

Clearly this is a better way of finding a hidden population than random surveys, since each named individual is likely to be a member of the population, and you also have the advantage of an introduction to them from one of their acquaintances, which may make it more likely that they will talk to you. However, there are some serious problems with the method as well. In particular, snowball sampling gives highly biased samples. In the limit of a large number of waves, snowball sampling samples actors non-uniformly with probability proportional to their “eigenvector centrality” (see Section 7.1.2). In principle, knowing this, one could compensate for the non-uniformity by appropriately weighting the results, but in practice the limit of large number of waves is rarely reached, and in any case the eigenvector centrality cannot be calculated without knowledge of the complete contact network, which by definition we don’t have, making correct weighting impossible. In short, snowball sampling gives biased samples of populations and there is little we can do about it. Nonetheless, the technique is sufficiently useful for finding populations that are otherwise hard to pin down that it has been widely used, biases and all, in studies over the past few decades.

Sometimes, in the case of small target populations, a few waves of snowball sampling may find essentially all members of a local population, in which case

the method can be regarded as returning data about the structure of the social network. If the contacts of each interviewed participant are recorded in the study, it should be possible to reconstruct the contact network when the study is complete. This has occasionally been done, although as noted above, the object is more often to exploit the social network to find the population than to study the network itself.

A technique closely related to snowball sampling is *contact tracing*, which is essentially a form of snowball sampling applied to disease incidence. Some diseases, such as tuberculosis and HIV, are considered in many countries to be sufficiently serious that, when someone is discovered to be carrying them, an effort must be made to track down all those who might also have been infected. Thus, when a patient tests positive for HIV, for instance, he or she will be questioned about recent sexual contacts, and possibly about other types of potentially disease-causing contacts, such as needle sharing if the patient is an injection drug user. Then health authorities will make an effort to track down the people so identified and test them for HIV also. The process is repeated with any who test positive, tracing their contacts as well, and so forth, until all leads have been exhausted. While the primary purpose of contact tracing is to curtail disease outbreaks and safeguard the health of the population, the process also produces data about the network through which a disease is spreading and such data have sometimes been used in scientific studies, particularly of sexually transmitted diseases, for which data may otherwise be hard to come by. Data from contact tracing studies display biases similar in type and magnitude to those seen in snowball sampling and should be treated with the same caution. Indeed, they may contain extra biases as well, since contacts are rarely pursued when an individual tests negative for the disease in question, so the sample is necessarily dominated by carriers of the disease, who are themselves usually a biased sample of the population at large.

There is another variant of snowball sampling that deals to some extent with the problems of bias in the sample. This is *random-walk sampling* [270,445]. In this method one again starts with a single member of the target community and interviews them to determine their contacts. Then, however, instead of tracking down all of those contacts, one chooses one of them at random and interviews only that one at the next step. If the person in question cannot be found or declines to be interviewed, one chooses another contact, and the process is repeated. Initially it appears that this will be a more laborious process than standard snowball sampling, since one spends a lot of time determining the names of individuals one never interviews, but this is not the case. In either method one has to determine the contacts of each person interviewed, so the total amount of work for a sample of a given size is the same. It is, however, very

important that one really does determine all the contacts of each individual, even though most of the time only one of them is pursued. This is because for the method to work correctly one must make a truly random choice among the complete set of contacts, for example by rolling a die (or some modern electronic version thereof).

The advantage of the random-walk sampling method is that, as shown in Section 6.14.3, the asymptotic sampling probability of nodes in a random walk is simply proportional to node degree (see Eq. (6.44)). What's more, the asymptotic regime in such studies is, unlike snowball sampling, reached quite quickly for relatively small sample sizes.⁸

Knowing this, and given that we determine degree (i.e., the number of contacts an individual has) as a part of the interview process, we can easily compensate for sampling bias by a suitable weighting of the results and make population estimates of quantities in a way that is, in theory at least, unbiased. In practice, many sources of bias remain, particularly those associated with participant subjectivity, inability to recall contacts, and non-participation of named contacts. Still, random-walk sampling is a significant improvement on standard snowball sampling. Its principal disadvantage is that it is relatively slow. Since the participants are interviewed serially, in a chain, rather than in parallel waves, a strict implementation of the method can take a long time to develop a large sample. One can get around this obstacle to some extent by running several short random walks in parallel instead of one long one, but the walks cannot be too short or they will not reach the asymptotic regime in which sampling is proportional to degree. Moreover, random-walk sampling is still not a good technique for probing the structure of the contact network itself. It returns a linear chain of contacts, a walk through the network, not a picture of the overall network structure. It can be used to estimate some network quantities, such as clustering coefficients (Section 7.3) or degree distributions (Section 10.3), but not normally the complete structure.

Another variant of random-walk sampling is used to deal with a different problem, that of enrolling study participants. In some cases it is considered unethical to get participants to name their contacts, particularly when the topic of a study is one of dubious legality, and permission to perform such studies may be withheld by the authorities. To circumvent this problem one can make use of *respondent-driven sampling* [421]. In this technique, participants are

⁸In snowball sampling the sample size grows exponentially with the number of sampling waves and hence one typically only performs a logarithmic number of waves, which is not enough for the sampling process to reach equilibrium. In random-walk sampling the sample size grows only linearly and one must perform a linear number of random walk steps.

usually paid to take part, and enrollment is achieved by handing out tickets to interviewees. Rather than asking people to name their contacts, the interviewees are simply invited to give the tickets to their acquaintances, and told that both they and the acquaintances will receive payment if the acquaintance brings the ticket to the investigator and agrees to participate in the study. In this way, no one is ever asked to name names and all participants have actively volunteered their participation. In the case where a single ticket is given to each participant, the method is roughly equivalent to random-walk sampling and should in theory give a less biased sample than snowball sampling for the same reasons. In practice, a new bias is introduced because the recipient of the ticket is not necessarily chosen at random from an individual's acquaintances. Also, tickets frequently get lost or their recipients decline to participate, remuneration notwithstanding, so one normally gives out more than one ticket to each participant, which complicates the sampling process and can introduce further biases [413]. Even so, it is believed that respondent-driven sampling provides superior population samples to snowball sampling, and it is the method of choice for studies in which one cannot ask people to name their contacts.

CHAPTER 5

BIOLOGICAL NETWORKS

A discussion of networks of interest in biology, including metabolic networks, neural networks, and food webs

NETWORKS appear in many branches of biology as a convenient way of representing patterns of interaction between biological elements. Molecular biologists, for example, use networks to represent the patterns of chemical reactions among chemicals in the cell, while neuroscientists use them to represent patterns of connections between brain cells. In this chapter we describe the commonest kinds of biological networks and discuss methods for determining their structure.

5.1 BIOCHEMICAL NETWORKS

Biochemical networks, representing the molecular-level patterns of interaction and control in the biological cell, have attracted a significant amount of attention in recent years. The best studied examples are metabolic networks, protein–protein interaction networks, and genetic regulatory networks.

5.1.1 METABOLIC NETWORKS

Metabolism is the chemical process by which cells break down food or nutrients into usable building blocks and then reassemble those building blocks to form the biological molecules the cell needs to live. Typically, this breakdown and

reassembly involves chains or *pathways*, sets of successive chemical reactions that convert initial inputs into useful end products by a series of steps. The complete set of all reactions in all pathways forms a *metabolic network*, in which the nodes are the chemicals produced and consumed by the reactions—known generically as *metabolites*—and the edges are the reactions. By convention the definition of a metabolite is limited to small molecules, meaning things like carbohydrates (such as sugars) and lipids (such as fats), as well as amino acids and nucleotides. Amino acids and nucleotides are themselves the building blocks for larger polymerized macromolecules such as DNA, RNA, and proteins, but the macromolecules are not considered metabolites—they are not produced by simple chemical reactions but by more complex molecular machinery within the cell, and hence are treated separately. (We discuss some of the mechanisms by which macromolecules are produced in Section 5.1.3.)

Although the fundamental purpose of metabolism is to turn food into useful biomolecules, one should be wary of thinking of it simply as an assembly line, even a very complicated one. Metabolism is not just a network of conveyor belts in which one reaction feeds another until the final products fall out the end; it is a dynamic process in which the concentrations of metabolites can change widely and rapidly, and the cell has mechanisms for turning on and off the production of particular metabolites or even entire portions of the network. Metabolism is a complex machine that reacts to conditions both within and outside the cell and generates a broad array of chemical responses. A primary reason for the high level of scientific interest in metabolic networks is their importance as a stepping stone on the path towards an understanding of this machinery.

In general, an individual chemical reaction in the cell involves the consumption of one or more metabolites, which are broken down or combined to produce one or more others. The metabolites consumed are called the substrates of the reaction, while those produced are called the products. Most metabolic reactions, however, do not occur spontaneously, or do so only at a very low rate, so the cell employs an array of chemical catalysts, or *enzymes*, to make reactions occur at a usable rate. Unlike metabolites, enzymes are mostly macromolecules, usually proteins but occasionally RNAs. And like all catalysts, enzymes are not consumed in the reactions they catalyze but they play an important role in metabolism nonetheless. Not only do they enable reactions that would otherwise be thermodynamically disfavored or too slow to be useful, but they also provide one of the mechanisms by which the cell controls its metabolism. By increasing or decreasing the concentration of the enzyme that catalyzes a particular reaction, the cell can turn that reaction on or off, or moderate its speed. Enzymes tend to be highly specific to the reac-

tions they catalyze, each one enabling only one or a small number of reactions. Thousands of enzymes are known and thousands more are no doubt waiting to be discovered, and this large array of highly specific catalysts allows for a fine degree of control over the processes of the cell.

The details of metabolic networks vary between different species of organisms but, among animals at least, large parts are common to all or most species. Many important pathways, cycles, or other subportions of metabolic networks are essentially unchanged across the entire animal kingdom. For this reason one often refers simply to “metabolism” without specifying a particular species of interest; with minor variations, observations made in one species often apply to others.

The most natural network representation of a set of metabolic reactions is as a bipartite network. We encountered bipartite networks previously in Section 3.3.2 on recommender networks and in Section 4.5 on affiliation networks. A bipartite network has two distinct types of nodes, with edges running only between nodes of unlike kinds. In the case of metabolic networks the two types of nodes represent metabolites and reactions, with edges joining each metabolite to the reactions in which it participates. In fact, a metabolic network is really a *directed* bipartite network, since some metabolites go into the reaction (the substrates) and some come out of it (the products). By placing arrows on the edges we can distinguish between the ingoing and outgoing metabolites.¹ An example is sketched in Fig. 5.1a.

This bipartite representation of a metabolic network does not include any way of representing enzymes, which, though not metabolites themselves, are still an important part of metabolism. Although it’s not often done, one can in principle incorporate the enzymes by introducing a third class of nodes to represent them, with edges connecting them to the reactions they catalyze. Since enzymes are neither consumed nor produced in reactions, these edges are undirected—running neither into nor out of the reactions they participate in. An example of such a network is sketched in Fig. 5.1b. Technically this is now a *tripartite network*, partly directed and partly undirected.²

Correct and potentially useful though they may be, however, neither of these representations of metabolic networks finds much use. The most common representations of metabolic networks make use of a “projection” of the network onto just one set of nodes, either the metabolites or the reactions, with the

Projections of bipartite networks and the associated loss of information are discussed further in Section 6.6.

¹The metabolic network is the only example of a directed bipartite network appearing in this book, and indeed the only naturally occurring example of such a network that we are aware of, although no doubt there are others to be discovered if one looks in the right place.

²Also the only such network in the book.

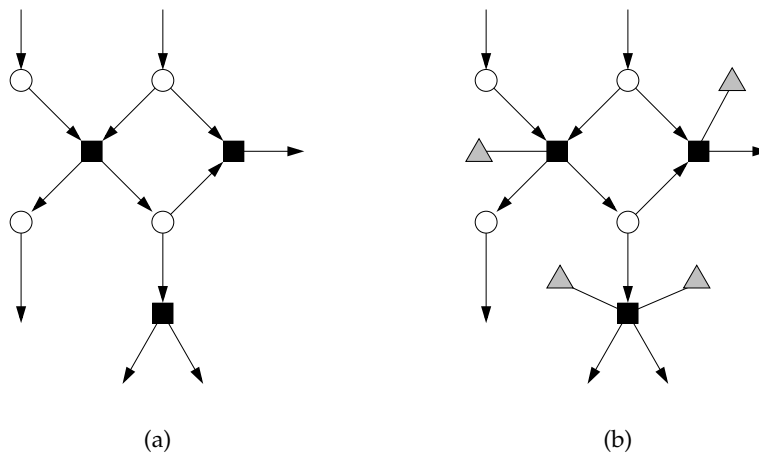


Figure 5.1: Bipartite and tripartite representations of a portion of a metabolic network. (a) A metabolic network can be represented as a directed bipartite network with nodes for the metabolites (circles) and reactions (squares) and directed edges indicating which metabolites are substrates (inputs) and products (outputs) of which reactions. (b) A third type of node (triangles) can be introduced to represent enzymes, with undirected edges linking them to the reactions they catalyze. The resulting network is a mixed directed/undirected tripartite network.

former being the more popular choice. In one approach the nodes in the network represent metabolites and there is an undirected edge between any two metabolites that participate in the same reaction, either as substrates or as products. Clearly this projection loses much of the information contained in the full bipartite network, though it is widely used nonetheless. Another approach, probably the most common, is to use a directed network with nodes representing the metabolites and a directed edge from one metabolite to another if there is a reaction in which the first metabolite appears as a substrate and the second as a product. This representation contains more of the information from the full network, but is still somewhat unsatisfactory since a reaction with many substrates or many products appears as many edges, with no easy way to tell that these edges represent aspects of the same reaction. The popularity of this representation arises from the fact that for many metabolic reactions only one product and one substrate are known or are considered important, and therefore the reaction can be represented by only a single directed edge with no confusion arising. A number of companies produce large charts showing

the most important parts of the metabolic network in this representation. An example is shown in Fig. 5.2. Such charts have become quite popular as wall decorations in the offices of molecular biologists and biochemists, although whether they are actually useful in practice is unclear.

The experimental measurement of metabolic networks is a complex and laborious process, although it has been made somewhat easier in recent years with the introduction of new techniques from molecular genetics. Experiments tend to focus neither on whole networks nor on individual reactions but on metabolic pathways. A number of tools are available to probe the details of individual pathways. Perhaps the most common is the use of radioactive isotopes to trace the intermediate products along a pathway. In this technique, the organism or cell under study is injected with a substrate for the pathway of interest in which one or more of the atoms has been replaced by a radioisotope. Typically, this has little or no effect on the metabolic chemistry, but as the reactions of the pathway proceed, the radioactive atoms move from metabolite to metabolite. Metabolites can then be refined, for example by mass spectroscopy or chromatography, and tested for radioactivity. Those that show it can be assumed to be downstream products in the pathway fed by the initial radioactive substrate.

This method tells us the products along a metabolic pathway, but of itself does not tell us in which order reactions take place along the pathway. Knowledge of the relevant biochemistry—which metabolites can be transformed into which others by chemical reactions—can help identify the ordering or at least narrow down the possibilities. Careful measurement of the strength of radioactivity of different metabolites, coupled with a knowledge of the half-life of the isotope used, can also give some information about pathway structure as well as rates of reactions.

Notice, however, that there is no way to tell if any of the reactions discovered have substrates other than those tagged with the radioisotope. If new substrates enter the pathway at intermediate steps (that is, they are not produced by earlier reactions in the pathway) they will not be radioactive and so will not be measured. Similarly, if there are reaction products that by chance do not contain the radioactive marker they too will not be measured.

An alternative approach to probing metabolic pathways is simply to increase the concentration in the cell of a substrate or enzyme for a particular reaction, thereby increasing the levels of the products of that reaction and those downstream of it in the relevant pathway or pathways, increases that can be measured to determine the constituents of the pathway. This technique has the advantage of being able to detect products other than those that carry a particular radioactive marker inherited from a substrate, but it is still inca-

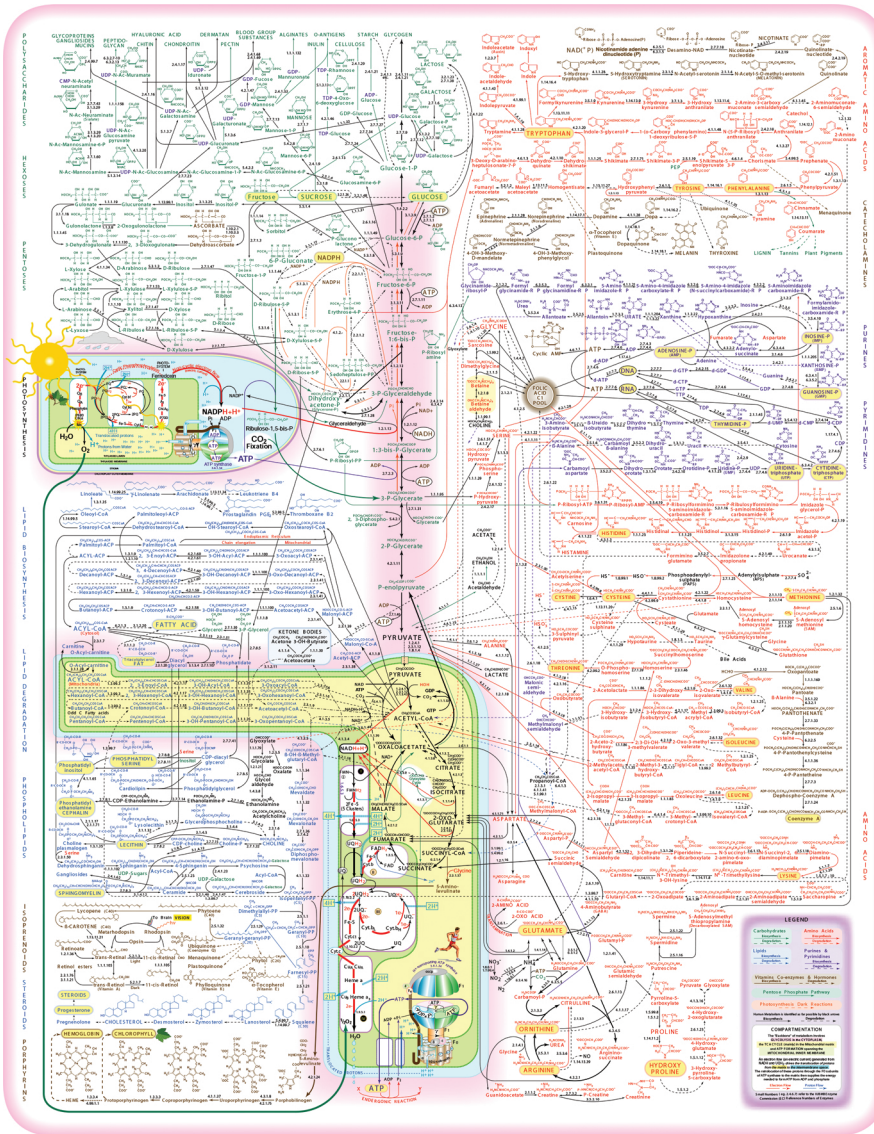


Figure 5.2: A metabolic network. A chart showing the network formed by the major metabolic pathways. Created by Donald Nicholson. Copyright of the International Union of Biochemistry and Molecular Biology. Reproduced with permission.

pable of identifying substrates other than those produced as products along the pathway.

A complementary experimental technique that can probe the substrates of reactions is reaction inhibition, in which a reaction in a pathway is prevented from taking place or its rate is reduced. Over time, this results in a build-up in the cell of the substrates for that reaction, since they are no longer being used up. By watching for this build-up one can identify the substrates. In principle the same method could also be used to determine the products of a reaction, since their concentration would decrease because they are not being produced any longer, but in practice this turns out to be a difficult measurement and is rarely done.

The inhibition of a reaction is usually achieved by disabling or removing an enzyme necessary for the reaction. This can be done in a couple of different ways. One can use *enzyme inhibitors*, which are chemicals that bind to an enzyme and prevent it from performing its normal function as a catalyst, or one can genetically alter the organism under study to remove or impair its ability to produce the enzyme (a so-called *knockout* experiment). The same techniques can also be used to determine which reactions are catalyzed by which enzymes in the first place, and hence to discover the structure of the third, enzymatic part of the tripartite metabolic network pictured in Fig. 5.1b.

The construction of a complete or partial picture of a metabolic network involves the combination of data from many different pathways, almost certainly derived from experiments performed by many different experimenters using many different techniques. There are now a number of public databases of metabolic pathway data which one can draw on to create networks, the best known being KEGG and MetaCyc. Assembling the network itself is a non-trivial task. Because the data are drawn from many sources, careful checking against the experimental literature is necessary to ensure consistent and reliable inputs to the process, and missing steps in metabolic pathways must often be filled in by guesswork based on biochemistry and a knowledge of the genetics. A number of computer software packages have been developed that can reconstruct networks from raw metabolic data in an automated fashion, but the quality of the networks they create is generally thought to be poorer than that of networks created by knowledgeable human scientists (although the computers are much faster).

5.1.2 PROTEIN-PROTEIN INTERACTION NETWORKS

The metabolic networks of the previous section describe the patterns of chemical reactions that turn one chemical into another in the cell. As we have

noted, the traditional definition of metabolism is restricted to small molecules and does not include proteins or other large molecules, except in the role of enzymes, in which they catalyze metabolic reactions but do not take part as reactants themselves.

Proteins do interact with one another and with other biomolecules, both large and small, but the interactions are not purely chemical. Proteins sometimes interact chemically with other molecules—exchanging small subgroups, for example, such as the exchange of a phosphate group in the process known as phosphorylation. But the primary mode of protein–protein interaction—interactions of proteins with other proteins—is physical, their complicated folded shapes interlocking to create so-called *protein complexes* (see Fig. 5.3) but without the exchange of particles or subunits that defines chemical reactions.

The set of all protein–protein interactions forms a *protein–protein interaction network*, in which the nodes are proteins and two nodes are connected by an undirected edge if the corresponding proteins interact. Although this representation of the network is the one commonly used, it omits much useful information about the interactions. Interactions that involve three or more proteins, for instance, are represented by multiple edges, and there is no way to tell from the network itself that such edges represent aspects of the same interaction. This problem could be addressed by adopting a bipartite representation of the network similar to the one we sketched for metabolic networks in Fig. 5.1, with two kinds of nodes representing proteins and interactions, and undirected edges connecting proteins to the interactions in which they participate. Such representations, however, are rarely used.

There are a number of experimental techniques available to probe for interactions between proteins. One of the most reliable and trusted is *co-immunoprecipitation*. Immunoprecipitation (without the “co-”) is a technique for extracting a single protein species from a sample containing more than one. The technique borrows from the immune system, which produces antibodies, specialized proteins that attach or *bind* to a specific other target protein when the two encounter each other. The immune system uses antibodies to neutralize proteins, complexes, or larger structures that are harmful to the body, but experimentalists have appropriated them for use in the laboratory. Immunoprecipitation involves attaching an antibody to a solid surface, such as the surface of a glass bead, then passing a solution containing the target protein (as well as others, usually) over the surface. The antibody and the target protein bind together, effectively attaching the protein to the surface via the antibody.

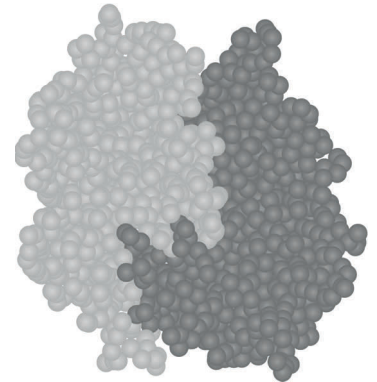
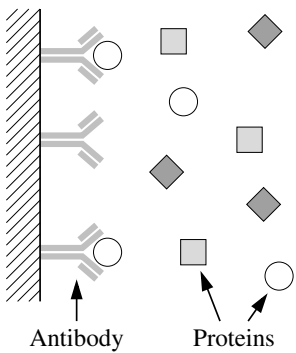


Figure 5.3: Two proteins joined to form a protein complex. Protein molecules can have complicated shapes that interlock with one another to form protein complexes.



In immunoprecipitation, antibodies attached to a solid surface bind to a specific protein, represented here by the circles, pulling it out of the solution.

Transcription factors are discussed in more detail in Section 5.1.3.

The rest of the solution washes away, leaving the target protein to be recovered from the surface.

There are known naturally occurring antibodies for many proteins of scientific interest, but researchers also routinely create antibodies for specific proteins by injecting those proteins (or more often a portion of a protein) into an animal to provoke its immune system to generate the appropriate antibody.

Co-immunoprecipitation is an extension of the immunoprecipitation method to the identification of protein interactions. An antibody is again attached to a suitable solid surface and binds to a known protein in a sample. If that protein is attached to others, forming a protein complex, then the entire complex will end up attached to the surface and will remain after the solution is washed away. Then the complex can be recovered from the surface and the different proteins that make it up individually identified, typically by testing to see if they bind to other known antibodies (a technique known as a *Western blot*).

Although well-established and reliable, co-immunoprecipitation is an impractical approach for reconstructing large interaction networks, since individual experiments, each taking days, must be performed for every interaction identified. If appropriate antibodies also have to be created, the process would take even longer; the creation of a single antibody involves weeks or months of work, and costs a considerable amount of money too. As a result, the large-scale study of protein–protein interaction networks did not really take off until the adoption in the 1990s and early 2000s of so-called *high-throughput* methods for discovering interactions, methods that can identify interactions quickly and in a semi-automated fashion.

The oldest and best established of the high-throughput methods for protein interactions is the *two-hybrid screen*, put forward by Fields and Song in 1989 [178].³ This method relies on the actions of a specialized protein known as a *transcription factor*, which, if present in a cell, turns on the production of another protein, referred to as a *reporter*. The presence of the reporter can be detected by the experimenter by any of a number of relatively simple means. The idea of the two-hybrid screen is to arrange things so that the transcription factor is created when two proteins of interest interact, thereby turning on the reporter, which tells us that the interaction has taken place.

The two-hybrid screen relies on the fact that transcription factors are typically composed of two distinct parts, a so-called *binding domain* and an *activation domain*. It turns out that most transcription factors do not require the binding and activation domains to be actually attached to one another for the transcrip-

³Also called a *yeast two-hybrid screen* or Y2HS for short, in recognition of the fact that the technique is usually implemented inside yeast cells, as discussed later.

tion factor to work. If they are merely in close enough proximity production of the reporter will be activated.

In a two-hybrid screen, a cell, usually a yeast cell, is persuaded to produce two proteins of interest, each with one of the domains of the transcription factor attached to it. This is done by introducing *plasmids* into the cell, fragments of DNA that code for the proteins and domains. Then, if the two proteins in question interact and form a complex, the two domains of the transcription factor will be brought together and, with luck, will activate production of the reporter.

In a typical two-hybrid experiment, the protein attached to the binding domain of the transcription factor is a known protein (called the *bait* protein) whose interactions the experimenter wants to probe. Plasmids coding for a large number of other proteins (called *prey*) attached to copies of the activation domain are created, resulting in a so-called *library* of possible interaction targets for the bait. The plasmids for the bait and the library of prey are then introduced into a culture of yeast cells, with the concentration of prey calibrated so that at most one prey plasmid enters each cell in most cases. Cells observed to produce the reporter are then assumed to contain plasmids coding for prey proteins that interact with the bait and the plasmids are recovered from those cells and analyzed to determine the proteins they correspond to.

The two-hybrid screen has two important advantages over older methods like co-immunoprecipitation. First, one can employ a large library of prey and hence test for interactions with many proteins in a single experiment, and second, the method is substantially cheaper and faster than co-immunoprecipitation per interaction detected. Where co-immunoprecipitation requires one to obtain or create antibodies for every protein tested, the two-hybrid screen requires only the creation of DNA plasmids and their later sequence analysis, both relatively simple operations for an experimenter armed with the machinery of modern genetic engineering.

One disadvantage of the two-hybrid screen is that the presence of the two domains of the transcription factor attached to the bait and prey proteins can get in the way of the proteins interacting with one another and prevent the formation of a protein complex, meaning that some legitimate protein–protein interactions will not take place under the conditions of the experiment.

The principal disadvantage of the method, however, is that it is simply unreliable. It produces high rates of both false positive results—apparent interactions between proteins that in fact do not interact—and false negative results—failure to detect true interactions. By some estimates the rate of false positives may be as high as 50%, meaning that fully half of all non-interacting proteins are wrongly reported as interacting. This has not stopped a number

See Section 5.1.3 for a discussion of DNA coding of proteins.

of researchers from performing analyses on the interaction networks reconstructed from two-hybrid screen data, but the results should be viewed with caution. It is certainly possible that many or even most of the conclusions of such studies are substantially inaccurate.

An alternative and more accurate class of methods for high-throughput detection of protein interactions are the *affinity purification* methods (also sometimes called *affinity precipitation* methods). These methods are in some ways similar to the co-immunoprecipitation method described previously, but avoid the need to develop antibodies for each protein probed. In an affinity purification method, a protein of interest is “tagged” by adding a portion of another protein to it, typically by introducing a plasmid that codes for the protein plus tag, in a manner similar to the introduction of transcription factor domains in the two-hybrid screen. Then the protein is given the opportunity to interact with a suitable library of other proteins and a solution containing the resulting protein complexes (if any) passed over a surface to which are attached antibodies that bind to the tag. As a result, the tag, the attached protein, and its interaction partners are bound to the surface while the rest of the solution is washed away. Then, as in co-immunoprecipitation, the resulting complex or complexes can be analyzed to determine the identities of the interaction partners.

The advantage of this method is that it requires only a single antibody that binds to a known tag, and the same tag–antibody pair can be used in different experiments to bind different proteins. Thus, as with the two-hybrid screen, one need only generate new plasmids for each experiment, which is relatively easy, as opposed to generating new antibodies, which is slow and difficult. Some implementations of the method have a reliability comparable to that of co-immunoprecipitation. Of particular note is the method known as *tandem affinity purification*, which combines two separate purification stages and generates correspondingly higher-quality results. Tandem affinity purification is the source for some of the most reliable current data for protein–protein interaction networks.

As with metabolic reactions, there are now substantial databases of protein interactions available online, such as BioGRID, STRING, and IntAct, and from these databases interaction networks can be constructed for analysis. An example is shown in Fig. 5.4.

5.1.3 GENETIC REGULATORY NETWORKS

As discussed in Section 5.1.1, the small molecules needed by biological organisms, such as sugars and fats, are manufactured in the cell by the chemical

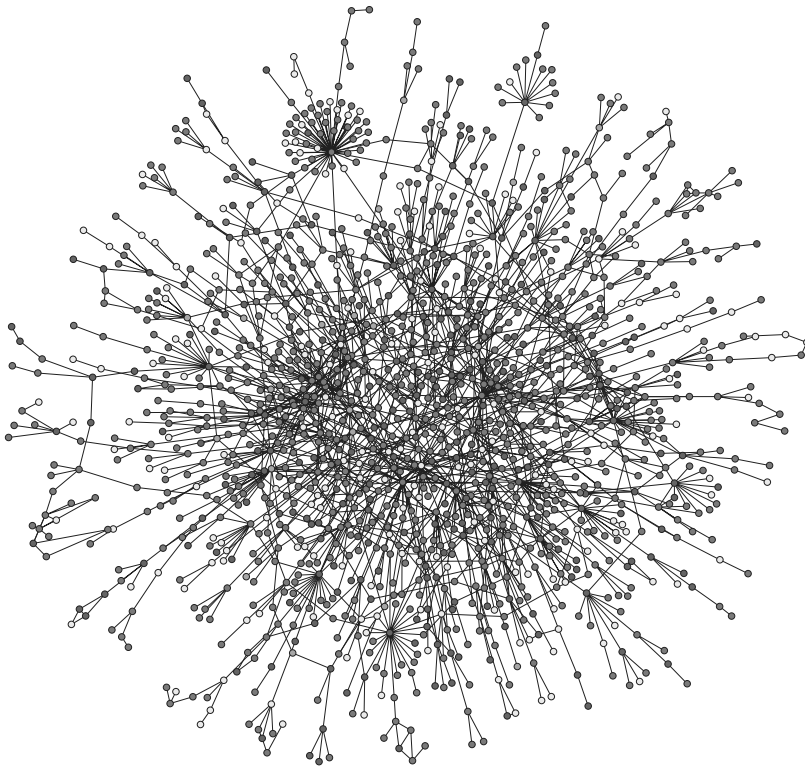


Figure 5.4: A protein–protein interaction network for yeast. A network of interactions between proteins in the single-celled organism *Saccharomyces cerevisiae* (baker’s yeast), as determined using, primarily, two-hybrid screen experiments. From Jeong *et al.* [250]. Copyright Macmillan Publishers Ltd., 2001. Reproduced by permission.

reactions of metabolism. Proteins, however, which are much larger molecules, are manufactured in a different manner, following recipes recorded in the cell’s genetic material, DNA.

Proteins are biological polymers, long-chain molecules formed by the concatenation of a series of basic units called *amino acids*. The individual amino acids themselves are manufactured by metabolic processes, but their assembly into complete proteins is accomplished by the machinery of genetics. There are 20 distinct amino acids that are used by all living organisms to build proteins, and different species of protein are distinguished from one another by the particular sequence of amino acids that make them up. Once created, a protein does not stay in a loose chain-like form, but folds up on itself under the

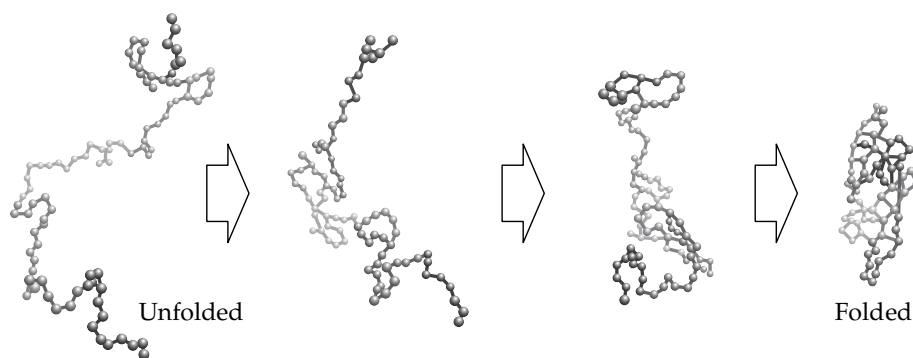


Figure 5.5: Protein folding. Proteins, which are long-chain polymers of amino acids, do not naturally remain in an open state (left), but collapse upon themselves to form a more compact folded state (right).

influence of thermodynamic forces and mechanical constraints, reliably producing a specific folded form or *conformation* whose detailed shape depends on the amino acid sequence—see Fig. 5.5. A protein’s conformation dictates the physical interactions it can have with other molecules and can expose particular chemical groups or active sites on the surface of the protein that contribute to its biological function within the organism.

A protein’s amino acid sequence is determined by a corresponding sequence stored in the DNA of the cell in which the protein is synthesized. This is the primary function of DNA in living matter, to act as an information storage medium containing the sequences of proteins needed by the cell. DNA is itself a long-chain polymer made up of units called *nucleotides*, of which there are four distinct kinds: adenine, cytosine, guanine, and thymine, commonly denoted A, C, G, and T, respectively.⁴ The amino acids in proteins are encoded in DNA as trios of consecutive nucleotides called *codons*, such as ACG or TTT, and a succession of such codons spells out the complete sequence of amino acids in a protein. A single strand of DNA can code for many proteins—hundreds or thousands of them—and two special codons, called the start and stop codons, are used to signal the beginning and end within the larger DNA strand of the

⁴Technically, DNA is a double-stranded polymer, having two parallel chains of nucleotides forming the famous double helix shape. However, the two strands contain essentially the same sequence of nucleotides and so for our purposes the fact that there are two is not important (although it is very important in other circumstances, such as in the reproduction of a cell by cellular division and in the repair of damaged DNA).

sequence coding for a protein. The DNA code for a single protein, from start codon to stop codon, is called a *gene*.

Proteins are created in the cell by a mechanism that operates in two stages. In the first stage, known as *transcription*, an enzyme called *RNA polymerase* makes a copy of the coding sequence of a single gene. The copy is made of RNA, another information-bearing biopolymer, chemically similar but not identical to DNA. RNA copies of this type are called *messenger RNAs*. In the second stage, called *translation*, the protein is assembled, step by step, from the RNA sequence by an ingenious piece of molecular machinery known as a *ribosome*, a complex of interacting proteins and RNA. The end result is a protein, assembled following the exact prescription spelled out in the corresponding gene. In the jargon of molecular biology, one says that the gene has been *expressed*.

The cell does not, in general, need to produce at all times every possible protein for which it contains a gene. Individual proteins serve specific purposes, such as catalyzing metabolic reactions, and it is important for the cell to be able to respond to its environment and circumstances by turning on or off the production of individual proteins as required. It does this by the use of *transcription factors*, which are themselves proteins and whose job is to control the transcription process by which DNA sequences are copied to RNA.

Transcription is performed by the enzyme RNA polymerase, which works by attaching to a DNA strand and moving along it, copying nucleotides one by one. The RNA polymerase doesn't just attach spontaneously, however, but is aided by a transcription factor. Transcription factors are specific to particular genes or sets of genes and regulate transcription in a variety of ways, but most commonly by binding to a recognized sub-sequence in the DNA, called a *promoter region*, which is adjacent to the beginning of the gene. The binding of the transcription factor to the promoter region makes it thermodynamically favorable for the RNA polymerase to attach to the DNA at that point and start transcribing the gene. Thus the presence in the cell of the transcription factor for the gene turns on or enhances the expression of that gene. We encountered an example of a transcription factor previously in our discussion of the two-hybrid screen in Section 5.1.2.

There are also transcription factors that inhibit expression by binding to a DNA strand in such a way as to prevent RNA polymerase from attaching to the strand and hence prevent transcription and the production of the corresponding protein.

Now we come to our main point: being proteins, transcription factors are themselves produced by transcription from genes. Thus, the protein encoded in a given gene can act as a transcription factor promoting or inhibiting production of one or more other proteins, which themselves can act as transcription factors

for further proteins and so forth. The complete set of such interactions forms a *genetic regulatory network*. The nodes in this network are proteins or equivalently the genes that code for them and a directed edge from gene A to gene B indicates that A regulates the expression of B. A slightly more sophisticated representation of the network distinguishes between promoting and inhibiting transcription factors, giving the network two distinct types of edges.

As with the metabolic networks of Section 5.1.1, genetic regulatory networks are an integral part of the machinery of the cell, a complex molecular mechanism capable of regulating many aspects of cellular behavior and coordinating a wide range of responses to environment changes, both within and outside the cell. Our knowledge of this mechanism is currently substantially incomplete, but it is hoped that a detailed analysis of the form and function of regulatory networks will help us understand more fully how organisms function at the molecular level.

Experimental determination of the structure of genetic regulatory networks involves identifying transcription factors and the genes they regulate. The process has several steps. To begin with, one first confirms that a given candidate protein does bind to DNA roughly in the region of a gene of interest. The commonest technique for doing this is the *electrophoretic mobility shift assay*, in which one creates strands of DNA containing the sequence to be tested and mixes them in solution with the candidate protein. If the two indeed bind, then the combined DNA-protein complex can be detected by *gel electrophoresis*, a technique in which one measures the speed of migration of electrically charged molecules or complexes through an agarose or polyacrylamide gel in an imposed electric field. In the present case the binding of the DNA and protein hinders the motion of the resulting complex through the gel, measurably reducing its speed when compared with unbound DNA strands. Typically, one runs two experiments side by side, one with protein and one without, and compares the rate of migration to determine whether the protein binds to the DNA. One can also run parallel experiments using many different DNA sequences to test which (if any) bind to the protein.

An alternative though less sensitive technique for detecting binding is the *deoxyribonuclease footprinting assay*. Deoxyribonucleases (also called DNases for short) are enzymes that, upon encountering DNA strands, cut them into shorter strands. There are many different DNases, some of which cut DNA only in particular places according to the sequence of nucleotides, but the footprinting technique uses a relatively indiscriminate DNase that will cut DNA at any point. If, however, a protein binds to a DNA strand at a particular location it will often prevent the DNase from cutting the DNA at or close to that location. Footprinting makes use of this by mixing strands of DNA containing

the sequence to be tested with the DNase and observing the resulting mix of strand lengths after the DNase has cut the DNA samples into pieces. Repeating the experiment with the protein present will result in a different mix of strand lengths if the protein binds to the DNA and prevents it from being cut in certain places. The mix is usually measured again by gel electrophoresis (strands of different lengths move at different speeds under the influence of the electric field) and one again runs side-by-side gel experiments with and without the protein to look for the effects of binding.

Both the mobility shift and footprinting assays can tell us whether a protein binds somewhere on a given DNA sequence. To pin down exactly where it binds one typically must do some further work. For instance, one can create short strands of DNA, called *oligonucleotides*, containing possible sequences that the protein might bind to, and add them to the mix. If they bind to the protein then this will reduce the extent to which the longer DNAs bind and visibly affect the outcome of the experiment. By a combination of such experiments, along with computer-aided guesswork about which oligonucleotides are likely to work best, one can determine the precise sub-sequence to which a particular protein binds.

While these techniques can tell us the DNA sequence to which a protein binds, they cannot tell us which gene's promoter region that sequence belongs to (if any), whether the protein actually affects transcription of that gene, or, if it does, whether the transcription is promoted or inhibited. Further investigations are needed to address these issues.

Identification of the gene is typically done not by experiment but by computational means and requires a knowledge of the sequence of the DNA in the region where the protein binds. If we know the DNA sequence then we can search it for occurrences of the sub-sequence to which our protein binds, and then examine the vicinity to determine what gene or genes are there, looking for example for start and stop codons in the region and then recording the sequence of other codons that falls between them. Complete DNA sequences are now known for many organisms, including humans, as a result of sequencing experiments starting in the late 1990s, and the identification of genes is, as a result, a relatively straightforward task.

Finally, we need to establish whether our protein actually acts as a transcription factor, which can be done either computationally or experimentally. The computational approach involves determining whether the sub-sequence to which the protein binds is indeed a promoter region for the identified gene. (It is possible for a protein to bind near a gene but not act as a transcription factor because the point at which it binds has no effect on transcription.) This is a substantially harder task than simply identifying nearby genes. The structure

of promoter regions is, unfortunately, quite complex and varies widely, but computer algorithms have been developed that can identify them with some success.

Alternatively, one can perform an experiment to measure directly the concentration of the messenger RNA produced when the gene is transcribed. This can be achieved, for example, by using a *microarray* (colloquially known as a “DNA chip”), tiny dots of DNA strands attached in a grid-like array to a solid surface. RNA will bind to a dot if a part of its sequence matches the sequence of the dot’s DNA and this binding can be measured using a fluorescence technique. By observing the simultaneous changes in binding on all the dots of the microarray, one can determine with some accuracy the change in concentration of any specific RNA and hence quantify the effect of the transcription factor. This technique can also be used to determine whether a transcription factor is a promoter or an inhibitor, something that is currently not easy using computational methods.

As with metabolic pathways and protein–protein interactions, there now exist electronic databases of genes and transcription factors from which it is possible to assemble snapshots of genetic regulatory networks. Current data on gene regulation are substantially incomplete and hence so are our networks, but more data are being added to the databases all the time.

5.1.4 OTHER BIOCHEMICAL NETWORKS

Metabolic, protein interaction, and regulatory networks are the best studied examples of biochemical networks, but there are a number of others that have been studied to a lesser extent. One example of potential medical importance is the *drug interaction network*. Pharmaceutical drugs are the first choice of treatment for most of the ailments that affect us, and it is not uncommon for people to take more than one drug at the same time, if they have multiple ailments for instance, or in the case of “combination therapies” that treat a single disease with a cocktail of two or more drugs. In cases like these one must be wary of potential interactions between drugs. For example, taking one drug may interfere with the action of another, decreasing its efficacy. Or the combination of two drugs may produce side effects, possibly serious, even when each drug alone is harmless.

The medical profession has assembled impressive databases of drug interactions, in which one can look up a particular drug and get a list of other drugs that should not be taken at the same time. Such a database can be turned into a drug interaction network, in which the nodes represent drugs and there is an edge between two drugs if they have a deleterious interaction [53]. A

more sophisticated network representation might also include labels, weights, or strengths on the edges to indicate the type or severity of the interaction. Analysis of drug interaction networks could tell us about patterns or regularities in the way drugs interact, or even allow us to predict previously unknown interactions [226].

A variation on the drug interaction network is the *drug–target network*. Medications typically act by binding to, activating, removing, enhancing, or inhibiting some chemical target in the body, such as a chemical receptor or a particular protein. Knowing which target a drug affects can be helpful in understanding its function and therapeutic potential, as well as informing us about possible drug interactions, since interactions often arise when two drugs affect the same target. Drugs and their targets can be represented as a bipartite network with two sets of nodes, one representing the drugs and the other the targets, and an edge connecting each target to the drugs known to affect it [53, 240]. Networks of drugs, targets, and interactions have received relatively little attention within the quantitative networks literature so far, but offer a promising avenue for future work that could have a significant impact on health and medicine.

Another application of network ideas in medicine is the *disease network* of Barabási *et al.* [42], which represents human diseases that have a genetic component. In this network the nodes represent diseases and there is an edge between two diseases if the same gene has been implicated in both. The result is a network in which related diseases, such as different forms of cancer, cluster together. As with the drug interaction network, one can also construct a bipartite variant of the disease network, in which there are two sets of nodes representing diseases and genes, and an edge connecting any gene to the diseases in which it is involved.

Networks have also been used as a representation of the structure of biomolecules themselves. As shown in Fig. 5.5, protein chains naturally fold up on themselves to create compact structures. In so doing, certain links in the chain end up in close proximity with others and form chemical bonds with them. It is largely the presence of these bonds, tying together different parts of the protein, that stabilizes the folded structure and gives it the characteristic shape that allows it to perform its biochemical function. The pattern of chemical bonds can be captured by a network in which the nodes represent amino acids—the links in the protein chain—and edges represent bonds between amino acid pairs, either the fundamental (or primary) bonds along the chain itself or the additional (secondary) bonds formed when the chain folds [83]. In practice, it may be difficult to tell exactly which amino acids interact with which others, in which case one can use simple spatial proximity as a proxy for interaction [35]. Similar networks can also be constructed for RNAs.

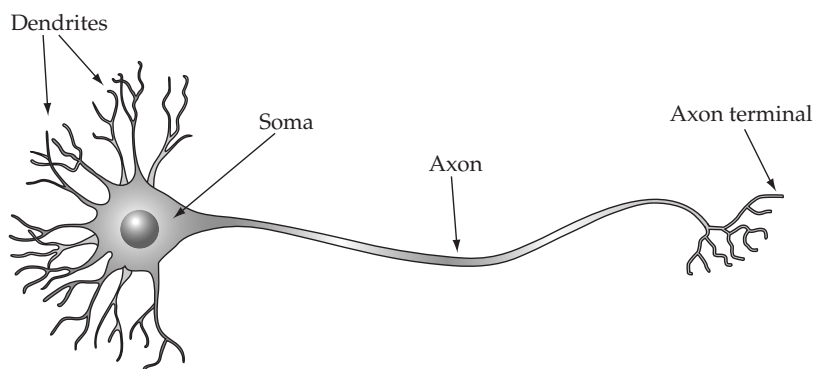


Figure 5.6: The structure of a neuron. A typical neuron is composed of a cell body or soma with many dendrites that act as inputs and a single axon that acts as an output. Towards its tip, the axon branches to allow it to connect to the inputs of several other neurons.

5.2 NETWORKS IN THE BRAIN

A completely different use of networks in biology arises in the study of the brain and central nervous system in animals. Two broad classes of brain networks are studied: microscopic networks of connections between individual brain cells and macroscopic networks of functional connection between entire brain regions.

5.2.1 NETWORKS OF NEURONS

One of the main functions of the brain is to process information, and the primary information processing element is the *neuron*, a specialized brain cell that combines (usually) several inputs to generate a single output. Depending on the species of animal, an entire brain can contain anywhere from a handful of neurons to more than a hundred billion, all wired together, the output of one cell feeding the input of another, to create a *neural network* capable of remarkable feats of calculation and decision making.

Figure 5.6 shows a sketch of a typical neuron, which consists of a cell body or *soma*, along with a number of protruding tentacles, which are essentially wires for carrying signals in and out of the cell. Most of the wires are inputs, called *dendrites*, of which a neuron may have just one or two, or as many as a thousand or more. Most neurons have only one main output, called the *axon*, which is typically longer than the dendrites and may in some cases extend over

large distances to connect the cell to others some way away. Although there is just one axon, it usually branches near its end to allow the output of the cell to feed the inputs of several others. The tip of each branch ends at an *axon terminal* that abuts the tip of the input dendrite of another neuron. There is a small gap, called a *synapse*, at the junction of the terminal and dendrite, across which the output signal of the first neuron must be conveyed in order to reach the second. The synapse plays an important role in the function of the brain, allowing the strength of the connection from one cell to another to be regulated by modifying the properties of the junction.⁵

The actual signals that travel within neurons are electrochemical in nature. They consist of traveling waves of electrical voltage created by the motion of positively charged sodium, calcium, or potassium ions in and out of the cell. These waves are called *action potentials* and involve voltage changes on the order of tens of millivolts traveling at tens of meters per second. When an action potential reaches a synapse, it cannot cross the gap between the axon terminal and the opposing dendrite by itself and the signal is instead transmitted chemically; the arrival of the action potential stimulates the release of a chemical neurotransmitter, which diffuses across the gap and activates receptor molecules at the other side. This in turn causes ions to move in and out of the dendrite, changing its voltage.

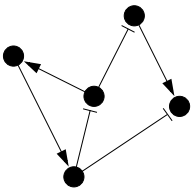
These voltage changes, however, do not yet give rise to another traveling wave. The soma of a neuron combines the inputs from its dendrites and as a result may (or may not) send an output signal down its own axon. The neuron is typically stable against perturbations caused by voltage changes at a small number of its inputs, but if enough inputs are excited they can collectively drive the neuron into an unstable runaway state in which it “fires,” generating a new action potential that travels down the cell’s axon, and so a signal is passed on to the next neuron or neurons in the network. Thus, the neuron acts as a switch or gate that aggregates the signals at its inputs and only fires when enough inputs are excited.

As described, inputs to neurons are excitatory, increasing the chance that the neuron fires, but inputs can also be inhibiting—signals received at inhibiting inputs make the receiving neuron less likely to fire. Excitatory and inhibiting inputs can be combined in a single neuron and the combination allows neurons to perform quite complicated information processing tasks all on their own, while

⁵Neurons do sometimes have direct connections between them without synapses. These direct connections are called *gap junctions*, a confusing name, since it sounds like it might be a description of a synapse but is in reality quite different. In our brief treatment of neural networks, however, we will ignore gap junctions.

an entire brain or brain region consisting of many neurons can perform tasks of extraordinary complexity. Current science cannot yet tell us exactly how the brain performs the more sophisticated cognitive tasks that allow animals to survive and thrive, but it is known that the brain constantly changes both the pattern and strength of the connections between neurons in response to inputs and experiences, and it is presumed that the details of these connections—the neural network—hold much of the secret. An understanding of the structure of neural networks is thus crucial if we are to explain the higher-level functions of the brain.

At the simplest level, a neuron can be thought of as a unit that accepts a number of inputs, either excitatory or inhibiting, combines them, and generates an output that is sent to one or more further neurons. In network terms, a neural network can thus be represented as a set of nodes—the neurons—connected by two types of directed edges, one for excitatory inputs and one for inhibiting inputs. (In this respect, neural networks are similar to the genetic regulatory networks of Section 5.1.3, which also contain both excitatory and inhibiting connections.) By convention, excitatory connections are denoted by an edge ending with an arrow “→”, while inhibiting connections are denoted by an edge ending with a bar “—|”.



A wiring diagram for a small neural network.

Neurons are not all the same. They come in a variety of different types and even relatively small regions or circuits in the brain may contain many types. This variation can be encoded in our network representation by different types of nodes. Visually the types are often denoted by using different shapes for the nodes or by labeling. In functional terms, neurons can differ in a variety of ways, including the number and type of their inputs and outputs, the nature and speed of their response to their inputs, whether and to what extent they can fire spontaneously without receiving inputs, and many other things besides.

Experimental determination of the structure of neural networks is difficult and the lack of straightforward experimental techniques for probing network structure is a major impediment to current progress in neuroscience. Some useful techniques do exist, however, although their application can be extremely laborious.

The basic tool for structure determination is microscopy, either optical or electronic. One relatively simple approach works with cultured neurons on flat dishes. Neurons taken from animal brains at an early stage of embryonic development can be cultivated in a suitable nutrient medium and will, without prompting, grow synaptic connections to form a network. If grown on a flat surface, the network is then roughly two-dimensional and its structure can be determined with reasonable reliability by simple optical microscopy. The advantage of this approach is that it is quick and inexpensive, but it has the

complete, depending on the size and complexity of the network studied.

Figure 5.7 shows an example of a “wiring diagram” of a neural network, reconstructed by hand from electron microscope studies of this type [470]. The network in question is the neural network of the worm *Caenorhabditis elegans*, one of the best studied organisms in biology. The brain of *C. elegans* is simple—it has less than 300 neurons and essentially every specimen of the worm has the same wiring pattern. Several types of neurons, denoted by shapes and labels, are shown in the figure, along with a number of different types of connections, both excitatory and inhibiting. Some of the edges run off the page, connecting to other parts of the network not shown. The experimenters determined the structure of the entire network and presented it as a set of interconnected wiring diagrams like this one.

The reconstruction of neural networks from slices in this way is the current gold standard in the field, but its laborious nature has led researchers to ask whether more direct methods of measurement might be possible. In the past few years a number of new methods have emerged that hold significant promise for faster and more accurate network structure determination. Most of these methods are based on optical (rather than electron) microscopy, which is something of a throwback to earlier days. Santiago Ramón y Cajal, the Nobel-prize-winning pathologist regarded by many as the father of neuroscience, pioneered the modern study of neuroanatomy with his beautiful hand-drawn illustrations of brain cells, created by staining slices of brain tissue with colored dyes and then examining them through an optical microscope (see Fig. 5.8). Current optical techniques do essentially the same thing, albeit with more technological sophistication.

Staining of brain tissue is crucial to making brain cells visible at optical wavelengths—without it there is not enough contrast between the neurons and surrounding tissue to make a clear picture. Early studies such as those of Ramón y Cajal used simple injected dyes, but modern studies use a range of more exotic techniques, particularly genetically engineered strains of laboratory animals, most often mice, that generate their own stains. This is usually done by introducing genes into the mice that produce fluorescent substances within brain cells such as the so-called *green fluorescent protein* or GFP, a widely used marker that was originally discovered, naturally occurring, in a certain species of jellyfish. Fluorescent proteins

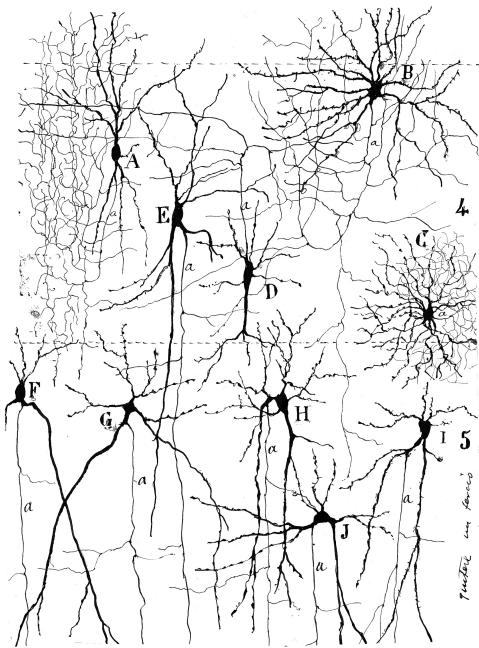


Figure 5.8: A historical neural network image. An early image of a collection of neurons, hand-drawn from optical microscope observations by Ramón y Cajal. Reproduced courtesy of the Cajal Institute: Cajal Legacy, Spanish National Research Council (CSIC), Madrid, Spain.

emit visible light when illuminated in the ultraviolet, light that can be photographed to create pictures of the neurons.

A crucial problem with optical imaging of brain tissue, however, is the sheer density of neurons; they are packed so tightly together—tens of thousands or more per cubic millimeter—that it is often difficult to tell them apart from one another. To get around this problem, researchers make use of a selection of different fluorescent proteins, including the original jellyfish GFP as well as various variants and alternatives, each emitting a different color of light. A particularly elegant implementation of this idea is the technique known as *Brainbow* [302], in which each neuron generates a random combination of different fluorescent proteins and each combination corresponds to a unique, identifiable color of emitted light. With a suitable palette of proteins the number of distinguishable colors can be as high as a hundred. The experimenter then makes separate images of the neurons with each of these colors, neurons which ideally are sparse enough to allow clear visualization of their shapes and positions, then combines the images to create a picture of the overall network.

While elegant, this approach does not solve the fundamental problem of having to slice up the brain to photograph it. *Brainbow* and techniques like it are still, at least for now, most often applied to slices. However, our new-found ability to clearly distinguish brain cells using optical techniques does open the door to the possibility of true 3D imaging if one can find a way to perform optical microscopy on whole brains or brain regions (something that is fundamentally impossible with electron microscopy). The fundamental tool for doing this is the *confocal microscope*, a type of microscope that uses special optics, combined with computer post-processing, to image the light coming from just a single two-dimensional slice of a three-dimensional space. By scanning the imaged slice through a sample one can then build up a picture of the entire three-dimensional structure. This doesn't completely solve our problem, however, because in order to focus light from a region in the interior of a brain the light still needs to get out of the brain in the first place, which normally it cannot do because the rest of the brain is in the way. One promising approach for resolving this issue is the technique called *Clarity* [105], which is a method for rendering brain tissue transparent by infusing it with a hydrogel. Once the tissue becomes transparent one can photograph its entire three-dimensional structure with a confocal microscope without needing to slice it up.

Methods such as these can allow one to visualize the positions and shapes of neurons in brains or brain regions, but they do not directly give us the topology of the corresponding neural network. For that, one must carefully analyze the pictures taken, following the path of each axon or dendrite to determine which neurons connect with which. And while this is certainly possible, it is

a laborious and sometimes error-prone task with current techniques. A quite different approach, which directly measures connections between neurons, is *transsynaptic tracing*, which involves injecting a tracer molecule—most commonly wheat germ agglutinin or WGA—into the brain, where it is absorbed by a subset of the neurons then transported along the axons of those neurons, across the synapses, and into the neighboring cells. In one ingenious version of the method the WGA is tagged with green fluorescent protein so that its final distribution can be photographed directly, from which one can then work out to which neighbors the outputs of a neuron connect. A variant on the same idea, called *retrograde tracing*, makes use of tracers that are naturally transported backwards across the synapse, allowing one to determine inputs. In more recent versions of these approaches researchers have replaced tracers like WGA with viruses that infect the neurons and spread from one to another, again allowing one to determine which cells are connected to which.

Optical imaging and transsynaptic tracing techniques are promising but still in their infancy. There is not yet (at the time of writing) any example of a large-scale network reconstruction, similar to that of Fig. 5.7, using these techniques. Still, this is a time of rapid advances in brain imaging and there is every hope that, probably within just a few years, these methods will have progressed to the point where they can give us significant insight into the structure of neural networks.

5.2.2 NETWORKS OF FUNCTIONAL CONNECTIVITY IN THE BRAIN

A different class of brain networks are networks of macroscopic functional connectivity between large-scale regions of the brain. In these networks the nodes represent entire brain regions, usually regions that are already known to perform some function such as vision, motor control, or learning and memory, and the edges represent some kind of functional connection, often only loosely defined, whereby one region controls or feeds information to another. The structure of these macroscopic networks can shed light on the logical organization of the brain—how information processing occurs or how different processes are interlinked—while avoiding the microscopic details of connection between individual brain cells. In principle macroscopic brain networks, while still complex, are much simpler than neuronal networks, the former containing typically tens or hundreds of nodes, where the latter could potentially contain billions. Macroscopic networks also have the advantage that they can be observed in living brains, including in humans, which cannot currently be done for their microscopic counterparts.

The primary technique for observing macroscopic network structure in the

brain is *magnetic resonance imaging*, or MRI, and particularly the techniques known as diffusion MRI and functional MRI. These are non-invasive imaging methods that can create a picture of a living brain from outside the skull. You may well have seen MRI images of the brain on TV or on the Internet, or even had an MRI done on your own brain. A fundamental disadvantage of MRI is its lack of spatial resolution: it can typically only resolve features on a scale of millimeters or larger, which is far bigger than the micrometer or nanometer scales of individual neurons. Nonetheless, for establishing the gross anatomical structure of the brain and its interconnection patterns at large scales, MRI is a useful tool.

Diffusion MRI (also called *diffusion tractography* or *diffusion-weighted MRI*, and sometimes abbreviated DW-MRI) aims to pick out physical connections between brain regions, the long-range wiring within the brain, which takes the form of bundles of elongated axons large enough to be resolved by MRI studies. Diffusion MRI detects anisotropies in the diffusion of water molecules. In elongated structures like axons, diffusion is faster along the axon than perpendicular to it, and the MRI can pick out this difference and hence determine the locations of the axons. Thus diffusion MRI can very directly pick out the edges in the macroscopic brain network, allowing us to reconstruct network topology.

By contrast, *functional MRI* (or fMRI) is a time-resolved imaging technique that picks out actual brain activity within living brain tissue in real time. Typically, it is sensitive to changes in blood oxygen level, which increases in active areas of the brain, causing them to “light up” on the MRI image. Functional MRI does not directly measure network structure the way diffusion MRI does. Instead one must infer connections by observing correlations between activity in different parts of the brain: two parts that routinely light up at the same time may be involved in the same tasks and hence are likely to be connected to one another. The combination of diffusion MRI and functional MRI together has, in recent years, allowed researchers to develop sophisticated maps of large-scale networks within both human and animal brains [86,436,454].

5.3 ECOLOGICAL NETWORKS

The final class of biological networks that we consider in this chapter is networks of ecological interaction between species. Species in an ecosystem can interact in a number of different ways. They can eat one another, they can parasitize one another, they can compete for resources, or they can have mutually advantageous interactions, such as pollination or seed dispersal. In principle, interactions of all of these types could be represented simultaneously in a

combined interaction network (perhaps a multilayer network—see Section 6.7). Traditionally, however, ecologists have separated interaction types into different networks. Networks of predator–prey interactions (i.e., who eats whom) have a particularly long history of study. Networks of hosts and parasites or of mutualistic interactions are less well studied, but have received some attention in recent years.

5.3.1 FOOD WEBS

The biological organisms on our planet can be divided into *ecosystems*, groups of organisms that interact with one another and with elements of their environment such as sources of material, nutrients, and energy. Mountains, valleys, lakes, islands, and larger regions of land or water can all be home to ecosystems composed of many organisms each. Within ecological theory, ecosystems are usually treated as self-contained units with no outside interactions, although in reality perfect isolation is rare and many ecosystems are only approximately self-contained. Nonetheless, the ecosystem concept is one of significant practical utility for understanding ecological dynamics.

A *food web* is a directed network that represents which species prey on which others in a given ecosystem.⁶ The nodes in the network correspond to species and the directed edges to predator–prey interactions. Figure 5.9 shows a small example, representing predation among species in Antarctica. There are several points worth noting about this figure. First, note that in this case not all of the nodes actually represent single species. Some of them do—the nodes for sperm whales and humans, for instance. But some of them represent collections of species, such as birds or fish. This is common practice in the study of food webs. If a set of species such as birds all prey upon and are preyed on by the same other species, then the network can be simplified by representing them as a single node, without losing any information about who preys on whom. Indeed, even in cases where a set of species only have mostly, but not exactly, the same predators and prey, we still sometimes group them if we feel the benefits of the resulting simplification are worth a small loss of accuracy. A set of species with the same or similar predators and prey is sometimes referred to as a *trophic species*.

⁶In common parlance, one refers to a *food chain*, meaning a chain of predator–prey relations between organisms starting with some lowly organism at the bottom of the chain, such as a microbe of some kind, and working all the way up to some ultimate predator at the top, such as a lion or a human being. Only a moment’s reflection, however, is enough to convince us that real ecosystems are not just single chains but complete networks of interactions.

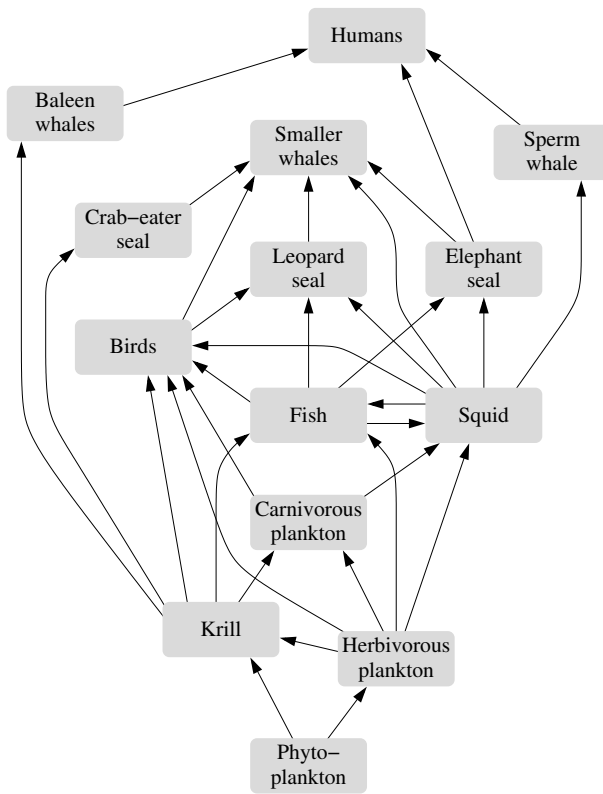


Figure 5.9: A food web of species in Antarctica. Nodes in a food web represent species or sometimes, as with some of the nodes in this diagram, groups of related species, such as fish or birds. Directed edges represent predator–prey interactions and run in the direction of energy flow, i.e., from prey to predator.

Second, note the direction of the edges in the network. One might imagine that the edges would point from predators to prey, but ecologists conventionally draw them in the opposite direction, from prey to predator. Thus the edge representing the eating of fish by birds runs *from* the fish node *to* the bird node. The reason for this apparently odd choice is that ecologists view food webs as representations of the flow of energy (or sometimes carbon) within ecosystems. The arrow from fish to birds indicates that the population of birds gains energy from the population of fish when the birds eat the fish.

Third, note that almost all the arrows in the figure run up the page. Directed networks with this property—that they can be drawn so that the edges all run

Acyclic networks are discussed in more detail in Section 6.4.1.

in one direction—are called *acyclic networks*. We encountered acyclic networks previously in our discussion of citation networks in Section 3.2. Food webs are usually only approximately acyclic. There are usually a few edges that do not run in the right direction,⁷ but it is often a useful approximation to assume that the network is acyclic.⁸

The acyclic nature of food webs indicates that there is an intrinsic pecking order among the species in ecosystems. Those higher up the order (which means higher up the page in Fig. 5.9) prey on those lower down, but not *vice versa*. A species' position in this pecking order is called by ecologists its *trophic level*. Species at the bottom of the food web, of which there is just one in our example—the phytoplankton—have trophic level 1. Those that prey on them—krill, herbivorous plankton—have trophic level 2, and so forth all the way up to the species at the top of the web, which have no predators at all. In our Antarctic example there are two species that have no predators, humans and small whales. (Note, however, that although such species are all, in a sense, at “the top of the food chain,” they need not have the same trophic level.)

Trophic level is a useful general guide to the roles that species play in ecosystems, those in lower trophic levels tending to be smaller, more abundant species that are prey to others higher up the food web, while those in higher trophic levels, the top predators, are usually larger-bodied and less numerous. Calculating a species' trophic level, however, is not always easy. In principle, the rule is simple: a species' trophic level is 1 greater than the trophic level of its prey. Thus, the herbivorous plankton and krill in our example have trophic level 2, because their prey has trophic level 1, and the carnivorous plankton have trophic level 3. But what happens if a species' prey do not all have the same trophic level? For instance, the squid in our example prey on species at two different levels, levels 2 and 3, so it is unclear what level the squid belong to. A variety of mathematical definitions have been proposed to resolve this issue. One strategy is to define trophic level to be 1 greater than the mean of the trophic levels of the prey. There is, however, no accepted standard definition, and the only indisputable statement one can make is that in most food webs some species have ill-defined or mixed trophic level.

⁷In Fig. 5.9, for example, there are edges in both directions between the fish and squid nodes, which makes it impossible to draw the network with all edges running in the same direction.

⁸Commonly omitted from food webs are the bacteria and other micro-organisms responsible for decomposing the bodies of dead animals and plants and feeding their nutrients back into the soil or the ocean bed. If one were to include these species, it would introduce additional energy flows from the top of the food web to the bottom and create closed loops, in which case the network would no longer be acyclic.

The food webs appearing in the ecological literature are of two main types. *Community food webs* are complete webs for an entire ecosystem, as in Fig. 5.9—they represent, at least in theory, every predator–prey interaction in the system. *Source food webs* and *sink food webs* are subsets of complete webs that focus on species connected, directly or indirectly, to a specific prey or predator. In a source food web, for instance, one records all species that derive energy from a particular source species, such as grass. Our food web of Antarctic species is, in fact, both a community food web and a source food web, since all of the species in the network derive their energy ultimately from phytoplankton. Phytoplankton is the source in this example, and the species above it (all of the species in this case) form the corresponding source web. A sink food web is the equivalent construct for a top predator in the network. In the Antarctic example, for instance, humans consume the sperm and baleen whales and elephant seals, which in turn derive their energy from fish, squid, plankton, krill, and ultimately phytoplankton. This subset of species, therefore, constitutes the sink food web for humans—the web that specifies through which species the energy consumed by humans passes.

The experimental determination of the structure of food webs is typically done in one of two different ways, or sometimes a mixture of both. The first and most straightforward method is direct measurement. Having settled on the ecosystem to be studied, one first assembles a list of the species in that ecosystem and then determines their predator–prey interactions. For large-bodied animals such as mammals, birds, or larger fish, some predation can be established simply by observation in the field—we see a bird eating a fish and the presence of the corresponding edge is thereby established. More often, however, and particularly with smaller-bodied animals, interactions are established by catching and dissecting the animals in question and examining the contents of their gut to determine what they have been eating.

The second primary method of constructing food webs is by compilation from existing literature. Many predator–prey interactions are already known and have been recorded in the scientific literature, but not in the context of the larger food web, and one can often reconstruct a complete or partial picture of a food web by searching the literature for such records. Many of the currently available food web data sets were assembled in this way, and some others were assembled by a combination of experimental measurement and literature searches. An interesting special case of the use of published records is in the construction of paleontological food webs. It turns out that some of the best-documented food webs today are not, in fact, from present-day ecosystems but from ecosystems that have been dead for millions of years, their food webs being assembled from the results of careful studies of fossilized species [152].

In some cases attempts have been made to measure not merely the presence (or absence) of interactions between species in ecosystems but also the strength of those interactions. One can quantify interaction strength by the fraction of its energy a species derives from each of its prey, or by the total rate of energy flow between a prey species and a predator. The result is a weighted directed network that sheds considerably more light on the flow of energy through an ecosystem than the more conventional unweighted food web. Measurements of interaction strength are, however, time-consuming and difficult, and yield uncertain results, so the data on weighted food webs should be treated with caution.

Food web data from a variety of sources have been assembled into publicly available databases such as Ecoweb [112], starting in the late 1980s.

5.3.2 OTHER ECOLOGICAL NETWORKS

Two other types of ecological network have received attention in the scientific literature (although less than has been paid to food webs). *Host–parasite networks* are networks of parasitic relationships between organisms, such as the relationship between a large-bodied animal and the insects and microorganisms that live on it and inside it. In a sense, parasitic relationships are a form of predation—one species eating another—but in practical terms they are quite distinct from traditional predator–prey interactions. Parasites, for example, tend to be smaller-bodied than their hosts where predators tend to be larger, and parasites can live off their hosts for long, sometimes indefinite, periods of time without killing them, whereas predation usually results in the death of the prey.

Parasitic interactions, however, do form networks that are somewhat similar to traditional food webs. Parasites themselves frequently play host to smaller parasites (called “hyperparasites”), which may have their own still smaller ones, and so forth through several levels.⁹ There is a modest but growing literature on host–parasite networks, much of it based on research within the agriculture community, a primary reason for interest in parasites being their prevalence in and effects on livestock and crop species.

The other main class of ecological networks is that of *mutualistic networks*,

⁹One is reminded of the schoolhouse rhyme by Augustus de Morgan:

Great fleas have little fleas upon their backs to bite 'em,
And little fleas have lesser fleas, and so ad infinitum.

meaning networks of mutually beneficial interactions between species. Three specific types of mutualistic network that have received attention in the ecological literature are networks of plants and the animals that pollinate them (primarily insects), networks of plants and the animals that disperse their seeds (such as birds), and networks of ant species and the plants that they protect and eat. Since the benefit of a mutualistic interaction runs, by definition, in both directions between a pair of species, mutualistic networks are undirected (or bidirectional, if you prefer), in contrast with the directed interactions of food webs and host–parasite networks. Most mutualistic networks (or at least most of those that have been studied) are also bipartite, consisting of two distinct, non-overlapping sets of species (such as plants and ants), with interactions only between members of different sets. In principle, however, non-bipartite mutualistic networks are also possible. The US National Center for Ecological Analysis and Synthesis has assembled an large collection of mutualistic network data in their Interaction Web Database, which can be found at <http://www.nceas.ucsb.edu/interactionweb>.

See Section 6.6 for a discussion of bipartite networks.

PART II

FUNDAMENTALS OF
NETWORK THEORY

CHAPTER 6

MATHEMATICS OF NETWORKS

An introduction to the mathematical tools used in the study of networks, tools that will be important to many subsequent developments

IN THE next several chapters we introduce the fundamental quantitative foundations of the study of networks, concepts that are crucial for the later developments in this book. In this chapter we discuss the basic theoretical tools used to describe and analyze networks, most of which come originally from graph theory, the branch of mathematics that deals with networks. Graph theory is a large field containing many results and we describe only a small fraction of them here, focusing on those most relevant to our goals in this book. Readers interested in pursuing the study of graph theory further might like to look at the books by Harary [230] or West [467].

In the three chapters after this one we look at measures and metrics for quantifying network structure (Chapter 7), computer algorithms for analyzing network data (Chapter 8), and statistical methods for networks (Chapter 9). Then in Chapter 10 we look at some of the remarkable patterns revealed in real-world networks when we apply the methods and metrics we have developed to their analysis.

6.1 NETWORKS AND THEIR REPRESENTATION

To begin at the beginning, a *network*—also called a *graph* in the mathematical literature—is, as we have said, a collection of nodes (or vertices) joined by

Networks, 2nd edition. Mark Newman, Oxford University Press (2018). © Mark Newman.
DOI: 10.1093/oso/9780198805090.001.0001

Network	Node	Edge
Internet	Computer or router	Cable or wireless data connection
World Wide Web	Web page	Hyperlink
Citation network	Article, patent, or legal case	Citation
Power grid	Generating station or substation	Transmission line
Friendship network	Person	Friendship
Metabolic network	Metabolite	Metabolic reaction
Neural network	Neuron	Synapse
Food web	Species	Predation

Table 6.1: Some examples of nodes and edges in networks

edges. Nodes and edges are also called *sites* and *bonds* in physics, and *actors* and *ties* in sociology.¹ Edges may also be variously called links, connections, or interactions, depending on context. Table 6.1 gives some examples of nodes and edges in particular networks.

Throughout this book we will normally denote the number of nodes in a network by n and the number of edges by m , which is a common notation in the mathematical literature. Most of the networks we will study have at most a single edge between any pair of nodes. In the rare cases where there can be more than one edge between the same nodes we refer to those edges collectively as a *multiedge*. In most networks there are also no edges that connect nodes to themselves, although such edges can occur in a few situations. Edges that connect nodes to themselves are called *self-edges* or *self-loops*.

A network that has neither self-edges nor multiedges is called a *simple network* or *simple graph*. A network with multiedges is called a *multigraph*. Figure 6.1 shows examples of (a) a simple graph and (b) a non-simple graph having both multiedges and self-edges.

There does not seem to be a special name given to networks with self-edges. They are just called “networks with self-edges.”

6.2 THE ADJACENCY MATRIX

The fundamental mathematical representation of a network is the *adjacency matrix*. Consider an undirected simple network with n nodes and let us label the nodes with integer labels $1 \dots n$, as we have, for instance, for the network in

¹This use of the word “actor” sometimes leads to confusion: an actor need not be a person who actually acts, and indeed need not be a person at all. In a social network of business relationships between companies, for instance, the actors are the companies (and the ties are the business relationships).

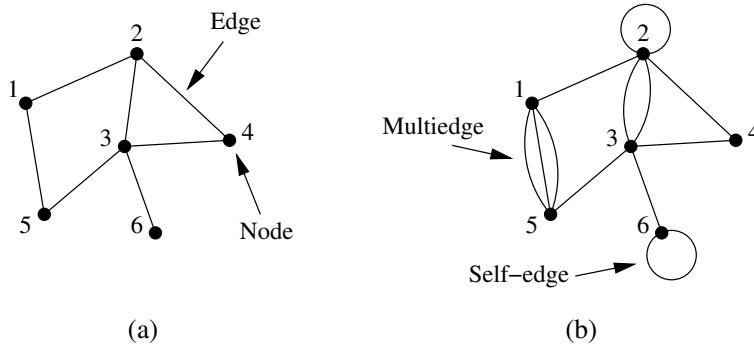


Figure 6.1: Two small networks. (a) A simple graph, i.e., one having no multiedges or self-edges. (b) A network with both multiedges and self-edges.

Fig. 6.1a. It does not matter which node gets which label, only that each label is unique, so that we can use the labels to refer to the nodes unambiguously.

The adjacency matrix \mathbf{A} of the network is now defined to be the $n \times n$ matrix with elements A_{ij} such that

$$A_{ij} = \begin{cases} 1 & \text{if there is an edge between nodes } i \text{ and } j, \\ 0 & \text{otherwise.} \end{cases} \quad (6.1)$$

For example, the adjacency matrix of the network in Fig. 6.1a is

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{pmatrix}. \quad (6.2)$$

Two points to note about the adjacency matrix are, first, that for a network such as this with no self-edges the diagonal matrix elements are all zero, and, second, that the matrix is symmetric, since if there is an edge between i and j then there is necessarily an edge between j and i .

It is also possible to represent multiedges and self-edges using an adjacency matrix. A multiedge is represented by setting the corresponding matrix element A_{ij} equal to the multiplicity of the edge. For example, a double edge between nodes i and j is represented by $A_{ij} = A_{ji} = 2$.

Self-edges are a little more complicated. A single self-edge from node i to itself is represented by setting the corresponding diagonal element A_{ii} of the

matrix equal to 2. Why 2 and not 1? Essentially, it is because a self-edge from i to i has two ends, both of which are connected to node i . As we will see, many mathematical results concerning the adjacency matrix work out more neatly if the matrix is defined this way, and thus it has become the accepted definition.² Another way to think about it is that non-self-edges appear twice in the adjacency matrix—an edge from i to j means that both A_{ij} and A_{ji} are 1. To count edges equally, self-edges should also appear twice, and since there is only one diagonal matrix element A_{ii} , we need to record both appearances there.

To give an example, the adjacency matrix for the multigraph in Fig. 6.1b is

$$\mathbf{A} = \begin{pmatrix} 0 & 1 & 0 & 0 & 3 & 0 \\ 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & 2 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 3 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 2 \end{pmatrix}. \quad (6.3)$$

One can also have multiple self-edges (or “multi-self-edges” perhaps). Such edges are represented by setting the corresponding diagonal element of the adjacency matrix equal to twice the multiplicity of the edge: $A_{ii} = 4$ for a double self-edge, 6 for a triple, and so forth.

6.3 WEIGHTED NETWORKS

Many of the networks we will study have edges that represent simple binary connections between nodes. Either they are there or they are not. In some situations, however, it is useful to represent edges as having a strength, weight, or value to them, usually a real number. Thus edges on the Internet might have weights representing the amount of data flowing along them or their bandwidth. In a food web, predator–prey interactions might have weights measuring total energy flow between prey and predator. In a social network connections might have weights representing frequency of contact between actors. Such *weighted* or *valued networks* can be represented mathematically by an adjacency matrix with the elements A_{ij} equal to the weights of the

²As discussed in Section 6.4, directed networks are different. In directed networks, self-edges are represented by a 1 in the corresponding diagonal element of the adjacency matrix.

corresponding connections. Thus, the adjacency matrix

$$\mathbf{A} = \begin{pmatrix} 0 & 2 & 1 \\ 2 & 0 & 0.5 \\ 1 & 0.5 & 0 \end{pmatrix} \quad (6.4)$$

represents a weighted network in which the connection between nodes 1 and 2 is twice as strong as that between 1 and 3, which in turn is twice as strong as that between 2 and 3.

Values on edges can also sometimes represent lengths of some kind. On a road or airline network, for instance, edge values could represent the number of kilometers or miles the edges cover, or they could represent travel time along the edges, which can be regarded as a kind of length—one denominated in units of time rather than distance. Edge lengths are, in a sense, the inverse of edge weights, since two nodes that are strongly connected can be regarded as “close” to one another and two that are weakly connected can be regarded as far apart. Thus one could perhaps convert lengths into weights by taking reciprocals and then use those values as elements of the adjacency matrix, although this should be regarded as only a rough translation; in most cases there is no formal mathematical relationship between edge weights and lengths.

We have now seen two different types of network where the adjacency matrix can have off-diagonal elements with values other than 0 and 1: networks with weighted edges and networks with multiedges. Indeed, if the weights in a weighted network are all integers it is possible to create a network with multiedges that has the exact same adjacency matrix, by choosing the multiplicities of the multiedges equal to the corresponding weights. This connection comes in handy sometimes. In some circumstances it is easier to reason about the behavior of a multigraph than a weighted network, or vice versa, and switching between the two can be a useful aid to analysis [355].

The weights in a weighted network are usually positive numbers, but there is no reason in theory why they could not be negative. Social networks of relations between people, for example, are sometimes constructed in which positive edge weights denote friendship or other cordial relationships and negative ones represent animosity.

And if edges can have weights on them, it is not a huge leap to consider weights on nodes too, or more exotic types of values on either edges or nodes, such as vectors or categorical variables like colors. Many such variations have been considered in the networks literature and we will discuss some of them later in the book. There is one other case of variables on edges, however, that is so central to the study of networks that we discuss it right away.

Networks with both positive and negative edges are discussed further in Section 7.5 when we consider the concept of structural balance.

6.4 DIRECTED NETWORKS

A *directed network* or *directed graph*, also called a *digraph* for short, is a network in which each edge has a direction, pointing *from* one node *to* another. Such edges are themselves called *directed edges*, or sometimes *arcs*, and can be represented graphically by, for instance, lines with arrows on them as in Fig. 6.2.

We have encountered a number of examples of directed networks in previous chapters, including the World Wide Web, in which hyperlinks run in one direction from one web page to another; food webs, in which energy flows from prey to predator; and citation networks, in which citations point from one paper to another.

The adjacency matrix of a directed network has matrix elements

$$A_{ij} = \begin{cases} 1 & \text{if there is an edge from } j \text{ to } i, \\ 0 & \text{otherwise.} \end{cases} \quad (6.5)$$

Note the direction of the edge here—it runs *from* the second index *to* the first. This is slightly counterintuitive, but it turns out to be convenient mathematically and it is the convention we adopt in this book.³

As an example, the adjacency matrix of the small network in Fig. 6.2 is

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix}. \quad (6.6)$$

Note that this matrix is not symmetric. In general the adjacency matrix of a directed network is asymmetric, since the existence of an edge from i to j does not necessarily imply that there is also an edge from j to i .

Like their undirected counterparts, directed networks can have multiedges and self-edges, which are represented in the adjacency matrix by elements with values greater than 1 and by non-zero diagonal elements, respectively. An important distinction, however, is that self-edges in a directed network are represented by setting the corresponding diagonal element to 1, not 2 as in the

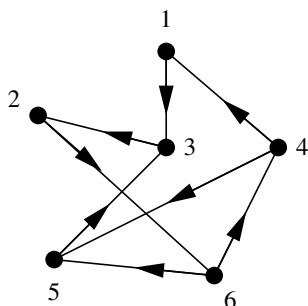


Figure 6.2: A directed network. A small directed network with arrows indicating the directions of the edges.

³Though common, this convention is not universal. One does sometimes see the opposite notation adopted, so one must be clear when reading (or writing) about directed networks which notation is in use.

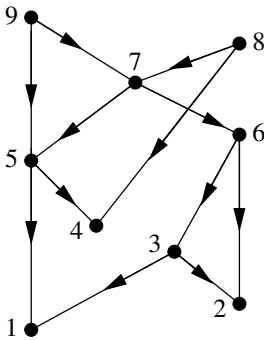


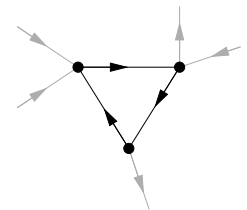
Figure 6.3: A directed acyclic network. In this network the nodes are laid out in such a way that all edges point downward. Networks that can be laid out in this way are called acyclic, since they possess no closed cycles of edges. A real-life example of an acyclic network is a network of citations between papers, in which the vertical axis would represent date of publication, running up the figure, and all citations would necessarily point from later papers to earlier ones.

undirected case. With this choice, formulas and results involving the adjacency matrix work out most neatly.

6.4.1 ACYCLIC NETWORKS

A *cycle* in a directed network is a closed loop of edges with the arrows on each of the edges pointing the same way around the loop. Networks like the World Wide Web have many such cycles in them. Some directed networks, however, have no cycles and these are called *acyclic* networks.^{4,5} A self-edge—an edge connecting a node to itself—counts as a cycle, so acyclic networks also have no self-edges.

A classic example of a directed acyclic network is a citation network of papers, as discussed in Section 3.2. When writing a paper you can only cite another paper if it has already been written, which means that all edges in a citation network point backwards in time, from later papers to earlier ones. Graphically we can depict such a network as in Fig. 6.3, with the nodes time-ordered—running from bottom to top of the picture in this case—so that all the edges point downward in the picture.⁶ There can be no closed cycles in such a network because any cycle would have to go down the picture and then come back up again to get back to where it started and there are no upward pointing edges with which to achieve this.



A cycle in a directed network.

⁴In the mathematical literature one often sees the abbreviation DAG, which is short for *directed acyclic graph*.

⁵Ones with cycles are called *cyclic*, although one doesn't often come across this usage, since directed networks are usually assumed to be cyclic unless otherwise stated.

⁶As discussed in Section 3.2, there are in real citation networks rare instances in which two papers both cite each other, forming a cycle of length two in the network, for instance if an author publishes two related papers in the same issue of a journal. Real citation networks are, thus, only approximately acyclic.

It is less obvious but still true that if a network is acyclic it can be drawn in the manner of Fig. 6.3 with all edges pointing downward. The proof that this can be done turns out to be useful, because it also provides us with a method for determining whether a given network is acyclic.

Suppose we have a directed acyclic network of n nodes. Then there must be at least one node somewhere in the network that has ingoing edges only and no outgoing ones. To see this suppose that it were not true and that every node has at least one outgoing edge. Then it would be possible to construct an endless path or “walk” through the network: we start at any node, follow one of the outgoing edges from that node, and repeat *ad infinitum*. But such a walk must, after at most n steps, revisit a node it has visited before (since there are only n nodes in total), and in so doing it completes a cycle in the network. But this cannot happen since our network is acyclic. Hence we have a contradiction and there must be at least one node with no outgoing edges.

Given this result, here is our scheme for drawing the network in ordered form, as in Fig. 6.3. First, we search through the network for a node with no outgoing edges. There could be more than one such node, in which case we choose whichever one we like. Let us call this node 1. We now remove node 1 from the network, along with any edges attached to it, then we repeat the process, finding another node with no outgoing edges in the remaining network. We call this node 2, remove it from the network along with its edges, and so forth.

After all nodes have been numbered and removed, we put the network back together again and draw a picture of it by placing the nodes in numerical order from bottom to top of the page and then drawing the directed edges in the appropriate positions between them. Every node has outgoing edges only to lower numbered nodes—those drawn below it in the picture—because it had no outgoing edges at the time it was removed from the network, so all its original outgoing edges (if it ever had any) must have been connected to nodes that were removed earlier. Thus all edges in the final picture must be pointing downward,⁷ achieving our goal of making a picture like that of Fig. 6.3.

This process is a useful one for visualizing acyclic networks. Most computer algorithms for drawing such networks work by arranging the nodes in order along an axis in just this way, and then moving them around along the other axis to make the network structure as clear and visually pleasing as possible

⁷Note that the particular order in which we draw the nodes, and hence the picture we produce, is not necessarily unique. If at any stage in the process there is more than one node with no outgoing edges then we have a choice about which one we remove and hence a choice between overall node orders.

(which usually means minimizing the number of times that edges cross).

The process is useful for another reason too: it will break down if the network contains cycles, and therefore it gives us a way to test whether a given network is acyclic. If a network contains a cycle, then none of the nodes in that cycle will ever be removed during our process: none of them will be without outgoing edges until at least one of the others in the cycle is removed, and hence none of them can ever be removed. Thus, if the network contains a cycle, there must come a point in our process where there are still nodes left in the network but all of them have outgoing edges. So a simple algorithm for determining whether a network is acyclic is:

1. Find a node with no outgoing edges.
2. If no such node exists, the network is *cyclic*. Otherwise, if such a node does exist, remove it and all its ingoing edges from the network.
3. If all nodes have been removed, the network is *acyclic*. Otherwise, go back to step 1.

The adjacency matrix of an acyclic directed network has interesting properties. Suppose we number the nodes of an acyclic network as described earlier, so that all edges point from higher numbered nodes to lower numbered ones. Then the adjacency matrix \mathbf{A} (whose element A_{ij} records the presence of an edge *from* j *to* i) has all its non-zero elements above the diagonal—it is upper triangular. For instance, the adjacency matrix of the network shown in Fig. 6.3 is

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}. \quad (6.7)$$

Note that the diagonal elements of the adjacency matrix are necessarily zero, since an acyclic network is not allowed to have self-edges. Triangular matrices with zeros on the diagonal are called *strictly triangular*.

If the nodes of an acyclic network are not numbered in the correct order as described earlier, then the adjacency matrix will not be triangular. (Imagine swapping rows and columns of the matrix above, for instance.) For every acyclic directed network, however, there exists at least one labeling of the nodes such that the adjacency matrix will be strictly upper triangular (and the algorithm described above can be used to find it).

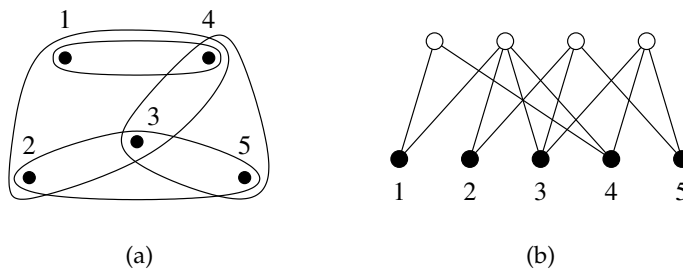


Figure 6.4: A hypergraph and corresponding bipartite graph. These two networks convey the same information—the membership of five nodes in four different groups. (a) The hypergraph representation in which the groups are represented as hyperedges, denoted by the loops circling sets of nodes. (b) The bipartite representation in which we introduce four new nodes (open circles at the top) representing the four groups, with edges connecting each of the original five nodes (bottom) to the groups to which it belongs.

6.5 HYPERGRAPHS

In some kinds of networks the edges join more than two nodes at a time. For example, we might want to create a social network representing family ties within a larger community of people. Families can have more than two people in them and one way to represent such families is to use a generalized type of edge that joins more than two nodes. Such an edge is called a *hyperedge* and a network with hyperedges is called a *hypergraph*.⁸ Figure 6.4a shows a small example of a hypergraph in which the hyperedges are denoted by loops.

Many of the networks that we will encounter in this book can be presented as hypergraphs. In particular, any network in which the nodes are connected together by common membership of groups of some kind, such as families, can be represented in this way. In sociology such networks are called “affiliation networks” and we saw several examples of them in Section 4.5. Directors sitting on the boards of companies, scientists coauthoring papers, and film actors appearing together in films are all examples of affiliation networks. See

⁸We could just use ordinary edges joining node pairs to represent our family ties, placing an edge between any two nodes that correspond to individuals in the same family. This, however, doesn’t tell us when two edges correspond to ties within the same family, and there is no single object in the network that corresponds to a family the way a hyperedge does in the hypergraph. In a number of ways, therefore, the hypergraph is a more complete representation of the pattern of family ties.

Network	Node	Group	Section
Film actors	Actor	Cast of a film	4.5
Coauthorship	Author	Authors of an article	4.5
Boards of directors	Director	Board of a company	4.5
Social events	People	Participants at social event	4.1
Recommender system	People	Those who like a book, film, etc.	3.3.2
Keyword index	Keywords	Pages where words appear	3.3.3
Rail connections	Stations	Train routes	2.4
Metabolic reactions	Metabolites	Participants in a reaction	5.1.1

Table 6.2: Hypergraphs and bipartite graphs. Examples of networks that can be represented as hypergraphs or equivalently as bipartite graphs. The last column gives the section of this book in which each is discussed.

Table 6.2 for more examples.

We will, however, talk very little about hypergraphs in this book, because there is another way of representing the same information that is more convenient for our purposes—the bipartite network.

6.6 BIPARTITE NETWORKS

A *bipartite network*, also called a *two-mode network* in the sociology literature, is a network with two kinds of nodes, and edges that run only between nodes of different kinds—see Fig. 6.5. Bipartite networks are most commonly used to represent the membership of a set of people or objects in groups of some kind. The people are represented by one set of nodes, the groups by the other, and edges join the people to the groups to which they belong. When used in this way a bipartite network captures exactly the same information as the hypergraphs of Section 6.5 (see Fig. 6.4) but for most purposes the bipartite graph is more convenient and it is certainly more widely used. We will use bipartite graphs frequently throughout this book.

For example, we can represent the network of film actors discussed in Section 4.5 as a bipartite network in which the two types of node are actors and films, and the edges connect actors to the films in which they appear. There are no edges that directly connect actors to other actors, or films to other films; the edges in a bipartite network only connect nodes of unlike kinds. As another example consider a recommender network, such as a network of who likes which books (see Section 3.3.2). The two types of nodes would then represent people and books, with edges connecting people to the books they like. Table 6.2 gives a number of further examples.

We discussed bipartite networks previously in the context of recommender networks in Section 3.3.2, affiliation networks in Section 4.5, and metabolic networks in Section 5.1.1.

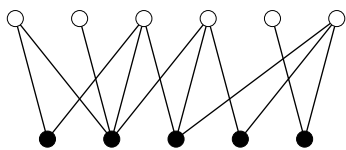


Figure 6.5: A small bipartite network. The open and closed circles represent two types of nodes and edges run only between nodes of different types. It is common to draw bipartite networks with the two sets of nodes arranged in lines, as here, to make the bipartite structure clearer. See Fig. 4.2 on page 50 for another example.

Bipartite networks do also occur occasionally in contexts other than membership of groups. For instance, there have been studies in the public health literature of networks of sexual contact—who sleeps with whom [271, 305, 392, 417]. If one were to construct such a network for a heterosexual population then the network would be bipartite, the two kinds of nodes corresponding to men and women and the edges corresponding to sexual contacts. (A network representing gay men or women on the other hand, or straight and gay combined, would probably not be bipartite.)

One occasionally also comes across bipartite networks that are directed. For example, the metabolic networks discussed in Section 5.1.1 can be represented as directed bipartite networks—see Fig. 5.1a. Weighted bipartite networks are also possible in principle, although no examples will come up in this book.

6.6.1 THE INCIDENCE MATRIX AND NETWORK PROJECTIONS

The equivalent of the adjacency matrix for an (undirected unweighted) bipartite network is a rectangular matrix called the *incidence matrix*. If n is the number of items or people in the network and g is the number of groups, then the incidence matrix \mathbf{B} is a $g \times n$ matrix having elements B_{ij} such that

$$B_{ij} = \begin{cases} 1 & \text{if item } j \text{ belongs to group } i, \\ 0 & \text{otherwise.} \end{cases} \quad (6.8)$$

For instance, the 4×5 incidence matrix of the network shown in Fig. 6.4b is

$$\mathbf{B} = \begin{pmatrix} 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 \end{pmatrix}. \quad (6.9)$$

Although a bipartite network may give the most complete representation of a particular system it is not always the most convenient. In some cases we would prefer to work with a network with only one type of node—a network of people alone, for instance, without the group nodes. One way to create such a network is to get rid of the group nodes and directly join together any two people who belong to the same group, creating a so-called *one-mode projection* of the two-mode bipartite form.

As an example, consider again the case of the films and actors. The one-mode projection onto the actors alone is the n -node network in which the nodes represent the actors and there is an undirected edge connecting any two actors who have appeared together in one or more films. We can also create a one-mode projection onto the films, which is the g -node network where the nodes represent films and two films are connected if they share one or more common actors. Every bipartite network has two one-mode projections in this way, one onto each of its types of nodes. Figure 6.6 shows the two one-mode projections of a small bipartite network.

When we form a one-mode projection, each group in the bipartite network results in a cluster of nodes in the projected network that are all connected to each other—a “clique” in network jargon (see Section 7.2.1). For instance, if a group in the bipartite network contains four members, then in the projection each of those four is connected to each of the others by virtue of common membership in the group. (Such a clique of four nodes is visible in the center of the lower projection in Fig. 6.6.) Thus, a one-mode projection is, generically, a union of a number of cliques, one for each group in the original bipartite network.

Projections are useful and widely employed, but their construction discards a lot of the information present in the original bipartite network and hence they are, in a sense, a less powerful representation of our data. For example, a projection loses any information about how many groups two nodes share in common. In the case of the actors and films, for instance, there are some pairs of actors who have appeared in many films together—Fred Astaire and Ginger Rogers, say, or William Shatner and Leonard Nimoy—and it’s reasonable to suppose this indicates a stronger connection than between actors who appeared together only once.

We can add information of this kind to our projection by making the projection weighted, giving each edge between two nodes in the projected network a weight equal to the number of common groups the nodes share. This weighted network still does not capture all the information in the bipartite original—it doesn’t record the total number of groups or the exact membership of each group for instance—but it is an improvement on the unweighted version and is quite widely used.

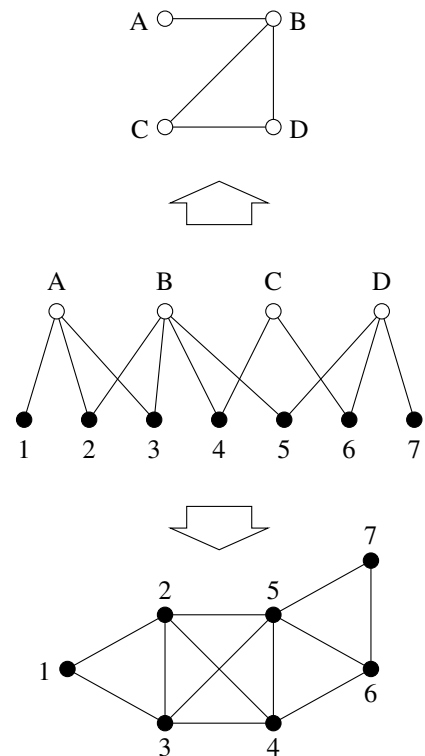


Figure 6.6: The two one-mode projections of a bipartite network. The central portion of this figure shows a bipartite network with four nodes of one type (open circles labeled A to D) and seven of another (filled circles, 1 to 7). At the top and bottom we show the one-mode projections of the network onto the two sets of nodes.

Mathematically, a one-mode projection can be written in terms of the incidence matrix \mathbf{B} of the original bipartite network as follows. The product $B_{ki}B_{kj}$ will be 1 if and only if i and j both belong to the same group k in the bipartite network. Thus, the total number P_{ij} of groups to which both i and j belong is

$$P_{ij} = \sum_{k=1}^g B_{ki}B_{kj} = \sum_{k=1}^g B_{ik}^T B_{kj}, \quad (6.10)$$

where B_{ik}^T is an element of the transpose \mathbf{B}^T of the incidence matrix. Equation (6.10) can be written in matrix notation as $\mathbf{P} = \mathbf{B}^T \mathbf{B}$, and the $n \times n$ matrix \mathbf{P} plays a role similar to an adjacency matrix for the weighted one-mode projection onto the n nodes. Its off-diagonal elements are equal to the weights in that network, the number of common groups shared by each node pair. \mathbf{P} is not quite an adjacency matrix, however, since its diagonal elements are non-zero, even though the one-mode network itself, by definition, has no self-edges. The diagonal elements have values

$$P_{ii} = \sum_{k=1}^g B_{ki}^2 = \sum_{k=1}^g B_{ki}, \quad (6.11)$$

where we have made use of the fact that B_{ki} only takes the values 0 and 1, so that $B_{ki}^2 = B_{ki}$. Thus P_{ii} is equal to the number of groups to which node i belongs.

To derive the adjacency matrix of the weighted one-mode projection, therefore, we would calculate the matrix $\mathbf{P} = \mathbf{B}^T \mathbf{B}$ and set the diagonal elements equal to zero. To derive the adjacency matrix of the unweighted projection, we would take the weighted matrix and replace every non-zero matrix element with a 1.

By a similar derivation, it is straightforward to show that the other one-mode projection, onto the groups, is represented by a $g \times g$ matrix $\mathbf{P}' = \mathbf{B} \mathbf{B}^T$, whose off-diagonal element P'_{ij} gives the number of common members of groups i and j , and whose diagonal element P'_{ii} gives the number of members of group i .

6.7 MULTILAYER AND DYNAMIC NETWORKS

The nodes and edges in a network need not all be of the same kind. Consider, for instance, a transportation network for a country or region in which nodes represent airports, train stations, bus stops, and so forth, while edges represent airline flights, train routes, etc. Such a structure could be captured by annotations on the nodes and edges describing their type, but a common and powerful alternative is to make use of a *multilayer network*.

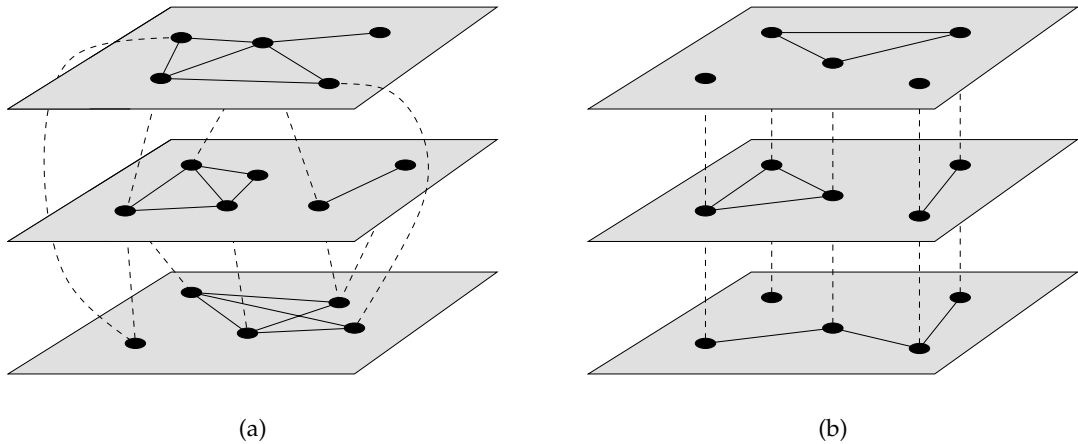


Figure 6.7: Multilayer and multiplex networks. (a) A multilayer network consists of a set of layers, each containing its own network, plus interlayer edges connecting nodes in different layers (dashed lines). An example is a transportation network with layers corresponding to airlines, trains, buses, and so forth. (b) A multiplex network is a special case of a multilayer network in which the nodes represent the same set of objects or people in each layer. For instance, a social network with several different types of connections could be represented as a multiplex network with one layer for each type. Dynamic or temporal networks are another example, where the layers represent snapshots over time of the structure of a single, time-varying network. In principle one can include interlayer edges in a multiplex network, as here, to represent the equivalence of nodes in different layers, although in practice these are often omitted.

A multilayer network is a set of individual networks, each representing nodes of one particular type and their connections, plus interlinking edges between networks—see Fig. 6.7a. The individual networks are referred to as layers. Thus, in our transportation example there might be a layer representing airports and flights, a layer representing train stations and train routes, and so on. Connections between layers could then be used to join nodes that are in the same geographical location, or at least close enough for easy walking. Many airports have train stations, for instance, and many train stations have bus stops. Paths through the resulting multilayer transportation network then represent possible passenger journeys: a passenger catches the bus to the train station for instance (represented by an edge in the bus route layer), walks from the bus stop into the station (an interlayer edge), then takes a train to their destination (an edge in the train route layer).

An important special case of a multilayer network occurs when the nodes in each layer represent the same set of points, objects, or individuals. Such networks, which are called *multiplex networks*, represent systems in which there

is only one type of node but more than one type of edge. An archetypal example is a social network, in which the nodes represent people, but there are many different kinds of connections between them—friendship connections, family connections, business connections, and so forth—each represented by a separate layer. The fact that the nodes represent the same people in every layer can be captured by interlayer edges connecting each node to its copies in other layers, although in practice such interlayer edges are often omitted for the sake of simplicity.

Another special case of a multilayer network is a *dynamic* or *temporal network*, a network whose structure changes over time. Most networks in fact do change over time—the Internet, the Web, social networks, neural networks, ecological networks, and many others change on a range of different time-scales. Most studies of networks ignore this fact and treat networks as static objects, which may be a reasonable approximation in some cases, but in others there is much to be learned by observing, analyzing, and modeling the time variation. Many empirical studies have been done on the way networks change over time. The usual approach is to measure the structure of the network repeatedly at distinct time intervals, resulting in a sequence of snapshots of the system, individual networks that can be thought of collectively as a multilayer network in which the layers have a specific ordering in time. In some examples only the edges change over time and not the nodes, in which case we have a multiplex network. In others the nodes can also appear or disappear, in which case a full multilayer network is needed to capture the structure, with different sets of nodes in different layers and interlayer edges between consecutive layers to indicate which nodes are equivalent. In some cases it may be useful to make the interlayer edges directed, pointing forward in time. For instance, when considering the spread of a disease over a time-varying contact network between individuals, possible routes the disease can take are represented by paths through the corresponding multilayer network, but the interlayer edges can be traversed only in the forward direction in time—catching the flu today can make you sick tomorrow but it cannot make you sick yesterday.

Mathematically, a multiplex network can be represented by a set of $n \times n$ adjacency matrices \mathbf{A}^α , one for each layer α (or each time point in the case of a dynamic network). Equivalently, one can think of the elements A_{ij}^α of these matrices as forming a three-dimensional tensor, and tensor analysis methods can usefully be applied to multiplex networks [264].

A multilayer network is more complicated. For a multilayer network, one must represent both the intralayer and interlayer edges, with potentially varying numbers of nodes in each layer. The intralayer edges can again be represented with a set of adjacency matrices \mathbf{A}^α , although the matrices need not all

be the same size now. If there are n_α nodes in layer α then the corresponding adjacency matrix has size $n_\alpha \times n_\alpha$. The interlayer edges can be represented by a set of additional *interlayer adjacency matrices*. The interlayer adjacency matrix $\mathbf{B}^{\alpha\beta}$ is an $n_\alpha \times n_\beta$ rectangular matrix with elements $B_{ij}^{\alpha\beta} = 1$ if there is an edge between node i in layer α and node j in layer β .

There are many examples of multilayer networks in empirical network studies. As we have said, most real-world networks are time-varying, and hence can be thought of as multiplex or dynamic networks [239]. Many social networks incorporate more than one type of interpersonal interaction and hence can be represented as multiplex networks. There have been a number of studies of transportation networks of the type discussed above, which are true multilayer networks [135, 198], and a range of other examples can be found in the literature. We refer the interested reader to the reviews by Boccaletti *et al.* [67], De Domenico *et al.* [134], and Kivela *et al.* [264], and the book by Bianconi [60].

6.8 TREES

A *tree* is a connected, undirected network that contains no loops—see Fig. 6.8a.⁹ By “connected” we mean that every node in the network is reachable from every other via some path through the network. A network can also consist of two or more parts, disconnected from one another, and if an individual part has no loops it is also called a tree. If all the parts of the network are trees, the complete network is called a *forest*.

Trees are often drawn in a *rooted* manner, as shown in Fig. 6.8b, with a *root node* at the top and a branching structure going down. The nodes at the bottom that are connected to only one other node are called *leaves*.¹⁰ Topologically, a tree has no particular root—the same tree can be drawn with any node, including a leaf, as the root node, but in some applications there are reasons for designating a specific root. A dendrogram is one example (see below).

Not many of the real-world networks that we encounter in this book are trees, although a few are. A river network is an example of a naturally occurring

The disconnected parts of a network are called “components”—see Section 6.12.

All trees are necessarily simple networks, with no multiedges or self-edges, since if they contained either then there would be loops in the network, which is not allowed.

⁹In principle, one could have directed trees as well, but the definition of a tree as a loopless network ignores edge directions if there are any. This means that a tree is not the same thing as a directed acyclic graph (Section 6.4.1), since the definition of a loop in a directed acyclic graph takes the directions of the edges into account. A directed acyclic graph may well have loops in it if we ignore directions (see, for example, Fig. 6.3).

¹⁰It may seem a little odd to draw a tree with the root at the top and the leaves at the bottom. Traditional trees of the wooden kind are, of course, the other way up. The upside-down orientation has, however, become conventional in mathematics and computer science, and we bow to that convention here.

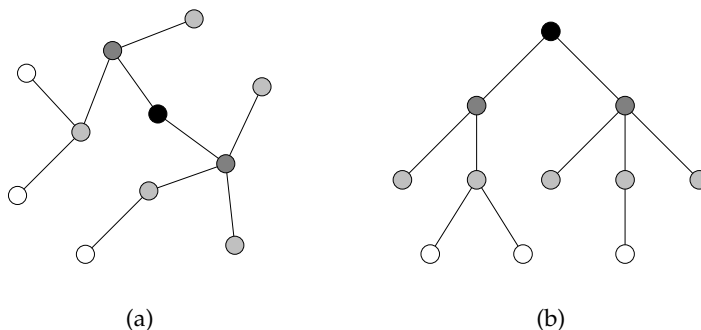


Figure 6.8: Two sketches of the same tree. The two panels here show two different depictions of a tree, a network with no closed loops. In (a) the nodes are positioned on the page in any convenient position. In (b) the tree is laid out in a “rooted” fashion, with a root node at the top and branches leading down to “leaves” at the bottom.

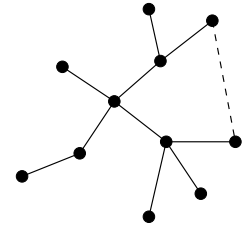
tree (see Fig. 2.6 on page 31). Trees do nonetheless play several important roles in the study of networks. In Chapter 11, for instance, we will study the network model known as the “random graph.” In this model local groups of nodes form trees and we can exploit this property to derive a variety of mathematical results about random graphs. In Section 14.5.1 we introduce the “dendrogram,” a useful tool that portrays a hierarchical decomposition of a network as a tree. Trees also occur commonly in computer science, where they are used as a basic building block for data structures such as AVL trees and heaps [9, 122] and in other theoretical contexts like minimum spanning trees [122], Cayley trees or Bethe lattices [388], and hierarchical models of networks (see Sections 14.7.2 and 18.3.2 and Refs. [109, 268, 465]).

Perhaps the most important property of trees for our purposes is that, since they have no closed loops, there is exactly one path between any pair of nodes. (In a forest there is at most one path, but there may be none.) This is clear since if there were two paths between a pair of nodes A and B then we could go from A to B along one path and back along the other, making a loop, which is forbidden.

This property of trees makes certain kinds of calculations particularly simple, and trees are sometimes used as a basic model of a network for this reason. For instance, the calculation of a network’s diameter (Section 6.11.1), the betweenness centrality of a node (Section 7.1.7), and certain other properties based on shortest paths are all relatively easy with a tree.

Another important property of trees is that a tree of n nodes always has exactly $n - 1$ edges. To see this, consider building up the tree by starting with a single node and no edges and adding further nodes one by one. With every node we must add at least one edge to keep the network connected, but on the other hand we cannot add more than one edge because if we did we would create a loop: the first edge we add connects the new node to the rest of the network but the second joins two nodes that are both already part of the network and hence already connected to one another. Adding an edge between two nodes that are already connected necessarily creates a loop (see figure), which is forbidden. Hence we must add exactly one edge to the network for every node we add. And since we start off with a single node and no edges it immediately follows that the tree always has one less edge and than it has nodes.

The reverse is also true, that any connected network with n nodes and $n - 1$ edges is a tree. If such a network were not a tree then there must be a loop in the network somewhere, implying that we could remove an edge without disconnecting any part of the network. Doing this repeatedly until no loops are left, we would end up with a tree, but one with less than $n - 1$ edges. As we showed above, however, every tree with n nodes must have $n - 1$ edges, and hence we have a contradiction. Thus we must have had a tree to begin with. As a corollary, this implies that the connected network on n nodes with the minimum number of edges is always a tree, since by the argument above no connected network has less than $n - 1$ edges and all networks with $n - 1$ edges are trees.



Adding an extra edge (dashed line) between any two nodes that are already part of the tree creates a loop.

6.9 PLANAR NETWORKS

A *planar network* is a network that can be drawn on a plane without having any edges cross.¹¹ Figure 6.9a shows a small planar network. Note that it is in most cases also possible to draw a planar network so that some edges do cross (Fig. 6.9b). The definition of planarity only specifies that at least one arrangement of the nodes exists that results in no crossing.

Most of the networks we will encounter in this book are not planar, but there are a few important examples that are. First of all, all trees are planar. For some trees, such as river networks, this is obvious. Rivers never cross one another; they only flow together. In other cases, such as the trees used in computer

¹¹A plane is a flat surface with open boundaries. One can define a generalization of a planar network for other types of two-dimensional surface, such as a torus, which wraps around on itself. A standard planar network, however, does not wrap around.

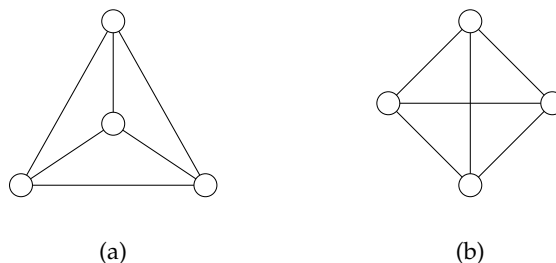


Figure 6.9: Two drawings of a planar network. (a) A small planar network with four nodes and six edges. It is self-evident that the network is planar, since in this depiction it has no edges that cross. (b) The same network redrawn with two of its edges crossing. Even though the edges cross, the network is still planar—a network is planar if it *can* be drawn without crossing edges. How you actually draw it is up to you.

data structures, there is no obvious two-dimensional surface onto which the network falls but it is planar nonetheless.

Among non-tree networks, some are planar for physical reasons. A good example is a road network. Because roads are confined to the Earth’s surface they form a roughly planar network. It does happen sometimes that roads meet without intersecting, one passing over another on a bridge, so that in fact, if one wishes to be precise, the road network is not planar. However, such instances are rare and the network is planar to a good approximation.

Another example is the network of which countries, states, or provinces are adjacent to which others—see Fig. 6.10. We can take a map depicting any set of contiguous geographic regions, represent each by a node, and draw an edge between any two that share a border. It is easy to see that the resulting network can always be drawn without crossing edges provided the regions in question are formed of contiguous landmasses.¹²

Networks of this type, representing regions on a map, have played an important role in mathematics, in the proof of the *four-color theorem*, which

¹²Technically, the map of the lower 48 US states in Fig. 6.10 does not quite satisfy this latter condition, since the state of Michigan is formed of two landmasses. (Several other states include offshore islands, but these are mostly too small to figure on our map.) We could get around this by having two nodes for the Upper and Lower Peninsulas of Michigan, though in that case we would no longer have exactly one node per state. In Fig. 6.10 we use just one node for Michigan, situated in the Lower Peninsula, but we do include an edge between Michigan and Wisconsin, which would not be present were it not for the Upper Peninsula, which shares a border with Wisconsin.

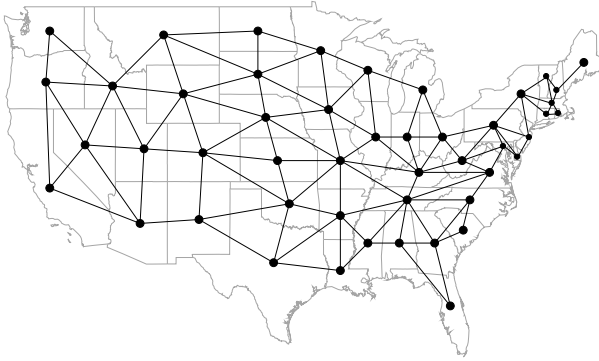


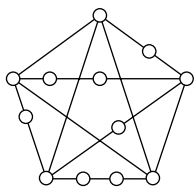
Figure 6.10: Graph of the adjacencies of the lower 48 United States. In this network each of the lower 48 states in the US is represented as a node and there is an edge between any two nodes if the corresponding states share a border. The resulting graph is planar, and indeed any set of states, countries, or other regions on a two-dimensional map can be turned into a planar graph in this way.

states that it is possible to color any set of regions on a two-dimensional map, real or imagined, with at most four colors such that no two adjacent regions have the same color, no matter how many regions there are or of what size or shape.¹³ By constructing the network corresponding to the map in question, this problem can be converted into a problem of coloring the nodes of a planar network in such a way that no two nodes connected by an edge have the same color. The number of colors required to color a network in this way is called the *chromatic number* of the network and many mathematical results are known about chromatic numbers. The proof of the four-color theorem—the proof that the chromatic number of a planar network is always four or less—is one of the triumphs of traditional graph theory and was first given by Appel and Haken [24–26] in 1976 after more than a hundred years of valiant effort within the mathematics community.¹⁴

An important question that arises in graph theory is how to determine, given

¹³The theorem only applies for a map on a surface with topological genus zero, such as a flat plane or a sphere. A map on a torus (which has genus 1) can require as many as seven colors.

¹⁴Appel and Haken’s proof was controversial at the time of its publication because it made extensive use of a computer to check large numbers of special cases. On the one hand, the proof was revolutionary for being the first proof of a major mathematical result generated in this fashion. On the other hand, a number of people questioned whether it could really be considered a proof at all, given that it was far too large for a human being to check its correctness by hand.



Technically, Kuratowski's theorem says that a non-planar network contains an *expansion* of K_5 or UG . An expansion is a network with any number (including zero) of extra nodes added along its edges, as in the expansion of K_5 shown here.

a particular network, whether that network is planar. For a small network it is a straightforward matter to draw a picture and play around with the positions of the nodes to see if one can find an arrangement in which no edges cross, but for a large network this is impractical and a more general method is needed. One such method makes use of *Kuratowski's theorem*, which states that every non-planar network must contain, somewhere within it, at least one of two distinctive smaller networks or subgraphs, called K_5 and UG , both of which are themselves non-planar. It immediately follows that a network is planar if, and only if, it contains neither of these subgraphs.

This approach is not, however, particularly useful for the analysis of real-world networks, because such networks are rarely precisely planar. (And if they are, then, as in the case of the shared border network of countries or states, it is usually clear for other reasons that they are planar and hence Kuratowski's theorem is unnecessary.) More often, like the road network, they are very nearly planar, but have a few edge crossings somewhere in the network. For such a network, Kuratowski's theorem would tell us, correctly, that the network was not planar, but we would be missing the point. What we would really like is some measure of the degree of planarity of a network, a measure that could tell us, for example, that the road network is 99% planar, even though there are a few bridges or tunnels here and there. One possible such measure is the minimum number of edge crossings with which the network can be drawn. This, however, would be a difficult quantity to calculate since, at least in the simplest approach, its evaluation would require us to try every possible way of drawing the network, of which there are an impossibly large number for all but the smallest of networks. Perhaps another approach would be to look at the number of occurrences of K_5 or UG in the network. So far, however, no widely accepted metric for degree of planarity has emerged. If such a measure were to gain currency it might well find occasional use in the study of real-world networks.

6.10 DEGREE

The *degree* of a node in an undirected network is the number of edges connected to it—see Fig. 6.11. In a social network of friendships between individuals, for instance, a person's degree is the number of friends they have. Note, however, that the definition of degree is in terms of number of edges, not number of neighboring nodes. The difference is important in multigraphs: if a node has two parallel edges to the same neighbor, both contribute to the degree—see Fig. 6.11b.

Despite its simplicity, degree is one of most useful and most widely used of

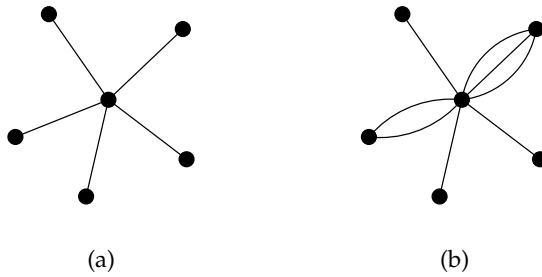


Figure 6.11: Degree of a node. (a) The central node has degree five because it has five attached edges. (b) The central node has five neighbors, but its degree is eight, because it has eight attached edges.

network concepts. It will play a large role in many of the developments in this book. Throughout the book we will denote the degree of node i by k_i . For a network of n nodes the degree can be written in terms of the adjacency matrix as¹⁵

$$k_i = \sum_{j=1}^n A_{ij}. \quad (6.12)$$

Every edge in an undirected network has two ends and if there are m edges in total then there are $2m$ ends of edges. But the number of ends of edges is also equal to the sum of the degrees of all the nodes, so

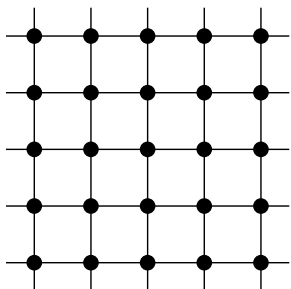
$$2m = \sum_{i=1}^n k_i = \sum_{ij} A_{ij}, \quad (6.13)$$

a result that we will use many times throughout this book.

The mean degree c of a node in an undirected network is

$$c = \frac{1}{n} \sum_{i=1}^n k_i, \quad (6.14)$$

¹⁵Note that this expression gives the correct result even if there are multiedges in the network, so long as the adjacency matrix is defined as in Section 6.2. It also works if there are self-edges, provided each self-edge edge is represented by a diagonal element $A_{ii} = 2$ as discussed in Section 6.2, and not 1.



An infinite square lattice is an example of a 4-regular network.

and combining this with Eq. (6.13) we get

$$c = \frac{2m}{n}. \tag{6.15}$$

This relation too will come up repeatedly throughout the book.

Occasionally we will come across networks in which all nodes have the same degree. In graph theory, such networks are called *regular graphs* or *regular networks*. A regular network in which all nodes have degree k is sometimes called *k-regular*. An example of a regular network is a periodic lattice such as a square or triangular lattice. On the square lattice, for instance, every node has degree four.

6.10.1 DENSITY AND SPARSITY

The maximum possible number of edges in a simple network (i.e., one with no multiedges or self-edges) is $\binom{n}{2} = \frac{1}{2}n(n-1)$. The *connectance* or *density* ρ of a network is the fraction of those edges that are actually present:

$$\rho = \frac{m}{\binom{n}{2}} = \frac{2m}{n(n-1)} = \frac{c}{n-1}, \tag{6.16}$$

where we have made use of Eq. (6.15) in the last equality. Most of the networks we are interested in are sufficiently large that (6.16) can be safely approximated as

$$\rho = \frac{c}{n}. \tag{6.17}$$

The density lies strictly in the range $0 \leq \rho \leq 1$. It can be thought of as the probability that a pair of nodes, picked uniformly at random from the whole network, is connected by an edge. This probability plays an important role in the random graph model discussed in Chapter 11.

Now consider a sequence of networks of increasing size n . If the density ρ remains non-zero as n becomes large the networks are said to be *dense*. In a dense network the fraction of non-zero elements in the adjacency matrix is non-vanishing in the limit of large n . A network where $\rho \rightarrow 0$ in the limit of large n is said to be *sparse*, and the fraction of non-zero elements in the adjacency matrix tends to zero.

These definitions only apply if you can actually take the limit $n \rightarrow \infty$, or at least extrapolate the limiting behavior from a sequence of networks of different sizes. When we are working with theoretical models of networks, as we will in later chapters of the book, we can take the limit formally and state whether a network is sparse or dense, but in practical situations involving observed

networks we cannot do this. We cannot take the limit as an empirical metabolic network or food web becomes large—we are stuck with the network nature gives us. For such networks there is no formal sense in which they are either sparse or dense.

Informally, on the other hand, one does often hear a network described, for example, as being sparse. Usually this just means that the value of ρ is small. In this qualitative sense, “sparse” just means that most of the possible edges that could exist in the network are not present.

In some cases real-world networks do change their sizes and by making measurements for different sizes we can make a guess as to whether they are best regarded as sparse or dense. The Internet and the World Wide Web are two examples of networks whose growth over time allows us to say with some conviction that they are best regarded as sparse.

In fact, most of the networks we examine in this book are usually considered to be sparse. There are very few examples where a network can truly be said to be dense, either in the mathematical sense above or in the more informal sense of just having a lot of edges.¹⁶ For our purposes, particularly when we come to study model networks, a more important distinction than that between sparse and dense is the distinction between networks with constant and diverging average degree.

Equation (6.17) tells us that the average degree c of a network is related to the density by $c = \rho n$, so in a dense network, where ρ is constant as $n \rightarrow \infty$, the average degree grows linearly with n , while for sparse networks the average degree grows sublinearly. And for some networks the average does not grow at all, meaning that ρ goes as $1/n$ for large n and c remains constant. Such networks will play an important role in the developments of this book. There seems to be no universally accepted name for them, although they are occasionally called *extremely sparse* [71].

Friendship networks, for example, plausibly have constant average degree, since it seems unlikely that the number of a person’s friends will be determined by the population of the world as a whole. How many friends a person has is more a function of how much time they have to devote to the maintenance of friendships, which is presumably independent of world population. Friendship networks therefore can be regarded as “extremely sparse.”

Arguably, indeed, most of the networks in this book fall into the extremely sparse category. If the average degree of a node does increase with n it usually

¹⁶A possible exception to the pattern is food webs. Studies comparing ecosystems of different sizes seem to show that the density of food webs is roughly constant, regardless of their size, indicating that food webs may be dense networks [153, 322].

does so only slowly, say as $\log n$. This sparsity has many implications. It makes possible a number of types of calculations that would otherwise be challenging, though at the same time it makes others harder. Sparsity will be particularly important when we look at computer algorithms in Chapter 8 and when we construct mathematical models of networks in Chapters 11 to 13.

6.10.2 DIRECTED NETWORKS

Node degrees are more complicated in directed networks. In a directed network each node has two degrees: the *in-degree* is the number of ingoing edges connected to a node and the *out-degree* is the number of outgoing edges. Bearing in mind that the adjacency matrix of a directed network has elements $A_{ij} = 1$ if there is an edge from j to i , the in- and out-degrees of node i can be written

$$k_i^{\text{in}} = \sum_{j=1}^n A_{ij}, \quad k_j^{\text{out}} = \sum_{i=1}^n A_{ij}. \quad (6.18)$$

These expressions also work for networks with multiedges, and for networks with self-edges provided a self-edge is represented by a diagonal element $A_{ii} = 1$ in the adjacency matrix, as discussed in Section 6.4.

The number of edges m in a directed network is equal to the total number of ingoing ends of edges at all nodes, or equivalently to the total number of outgoing ends of edges, so

$$m = \sum_{i=1}^n k_i^{\text{in}} = \sum_{j=1}^n k_j^{\text{out}} = \sum_{ij} A_{ij}. \quad (6.19)$$

Thus the mean in-degree c_{in} and the mean out-degree c_{out} of every directed network are equal:

$$c_{\text{in}} = \frac{1}{n} \sum_{i=1}^n k_i^{\text{in}} = \frac{1}{n} \sum_{j=1}^n k_j^{\text{out}} = c_{\text{out}}. \quad (6.20)$$

For simplicity we will just denote both by c and, combining Eqs. (6.19) and (6.20), we get

$$c = \frac{m}{n}. \quad (6.21)$$

Note that this differs by a factor of two from the equivalent result for an undirected network, Eq. (6.15).

6.11 WALKS AND PATHS

A *walk* in a network is any sequence of nodes such that every consecutive pair of nodes in the sequence is connected by an edge. In other words it is any route that runs from node to node along the edges. Walks can be defined for both directed and undirected networks. In a directed network, each edge traversed by a walk must be traversed in the direction of that edge. In an undirected network edges can be traversed in either direction.

In general a walk can intersect itself, revisiting a node it has visited before or running along an edge or set of edges more than once. Walks that do not intersect themselves are called *paths* or *self-avoiding walks*, and are important in many areas of network theory. Shortest paths and independent paths are two special cases of self-avoiding walks that we will study later.

The *length* of a walk in a network is the number of edges traversed along the walk (not the number of nodes). A given edge can be traversed more than once, and if so it is counted separately each time it is traversed. In layman's terms the length of a walk is the number of "hops" the walk makes from node to adjacent node.

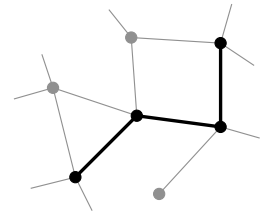
It is straightforward to calculate the number of walks of a given length r on a network. For either a directed or an undirected simple network the element A_{ij} is 1 if there is an edge from node j to node i , and 0 otherwise. (We will consider only simple networks for now, although the developments generalize easily to non-simple networks.) Then the product $A_{ik}A_{kj}$ is 1 if there is a walk of length 2 from j to i via k , and 0 otherwise. And the total number $N_{ij}^{(2)}$ of walks of length two from j to i , via any node, is

$$N_{ij}^{(2)} = \sum_{k=1}^n A_{ik}A_{kj} = [\mathbf{A}^2]_{ij}, \quad (6.22)$$

where $[\dots]_{ij}$ denotes the ij th element of the matrix.

Similarly the product $A_{ik}A_{kl}A_{lj}$ is 1 if there is a walk of length three from j to i via l and k , and 0 otherwise, and hence the total number of walks of length three is

$$N_{ij}^{(3)} = \sum_{k,l=1}^n A_{ik}A_{kl}A_{lj} = [\mathbf{A}^3]_{ij}. \quad (6.23)$$



A walk of length three in an undirected network.

Generalizing to walks of arbitrary length r , it is straightforward to see that¹⁷

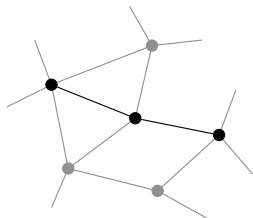
$$N_{ij}^{(r)} = [\mathbf{A}^r]_{ij}. \quad (6.24)$$

A special case of this result is that the number of walks of length r that start and end at the same node i is $[\mathbf{A}^r]_{ii}$. These walks are just loops in the network and the total number L_r of loops of length r in the network is the sum of this quantity over all possible starting points i :

$$L_r = \sum_{i=1}^n [\mathbf{A}^r]_{ii} = \text{Tr } \mathbf{A}^r. \quad (6.25)$$

Note that this expression counts separately loops consisting of the same nodes in the same order but with different starting points. Thus the loop $1 \rightarrow 2 \rightarrow 3 \rightarrow 1$ is considered different from the loop $2 \rightarrow 3 \rightarrow 1 \rightarrow 2$. The expression also counts separately loops that consist of the same nodes but traversed in opposite directions, so that $1 \rightarrow 2 \rightarrow 3 \rightarrow 1$ and $1 \rightarrow 3 \rightarrow 2 \rightarrow 1$ are distinct.¹⁸

6.11.1 SHORTEST PATHS



A shortest path of length two between two nodes.

A *shortest path* in a network, also sometimes called a *geodesic path*, is the shortest walk between a given pair of nodes, i.e., the walk that traverses the smallest number of edges. The *shortest distance* or *geodesic distance* between two nodes, often loosely called just the “distance,” is the length of the shortest path in terms of number of edges. In mathematical terms, the shortest distance between nodes i and j is the smallest value of r such that $[\mathbf{A}^r]_{ij} > 0$. (In practice, however, there are much better ways of calculating it than by employing this formula. We will study some of them in Section 8.5.)

Shortest paths and shortest distances play an important role in a number of network phenomena. For instance, the small-world effect discussed in Section 4.6 is in effect a statement about shortest distances—that they are surprisingly small even in the largest of networks. And shortest distances are

¹⁷For a more rigorous proof we can use induction. If there are $N_{ik}^{(r-1)}$ walks of length $r-1$ from k to i , then by arguments similar to those above there are $N_{ij}^{(r)} = \sum_k N_{ik}^{(r-1)} A_{kj}$ walks of length r from j to i , or in matrix notation $\mathbf{N}^{(r)} = \mathbf{N}^{(r-1)} \mathbf{A}$, where $\mathbf{N}^{(r)}$ is the matrix with elements $N_{ij}^{(r)}$. This implies that if $\mathbf{N}^{(r-1)} = \mathbf{A}^{r-1}$ then $\mathbf{N}^{(r)} = \mathbf{A}^r$. Starting from the base case $\mathbf{N}^{(1)} = \mathbf{A}$ we then have $\mathbf{N}^{(r)} = \mathbf{A}^r$ for all r by induction, and taking the ij th element of both sides gives Eq. (6.24).

¹⁸If we wish to count each loop only once, we should roughly speaking divide by r , but this does not allow for walks that have symmetries under a change of starting points, such as walks that consist of the same subloop traversed repeatedly. Counting such symmetric walks properly is a complex problem that can be solved exactly in only a few cases.

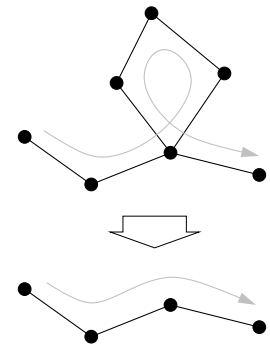
important in communication networks such as the Internet, where they affect how rapidly it is possible to get data from one node to another, or in transportation networks such as airline networks, where they determine how many legs will be required for a particular journey.

It is possible for there to be no shortest path between two nodes if the nodes are not connected together by any route through the network (i.e., if they are in different “components”—see Section 6.12). In this case one sometimes says that the distance between the nodes is infinite, although this is mostly just convention—it doesn’t really mean very much beyond the fact that the nodes are not connected.

Shortest paths are necessarily self-avoiding, which is why we call them paths. (Recall that a path means a self-avoiding walk.) If a walk intersects itself then it contains a loop and can be shortened by removing that loop while still connecting the same start and end points (see figure), and hence self-intersecting walks are never the shortest route between any two nodes.

Shortest paths are not necessarily unique, however. It is perfectly possible to have two or more paths of equal length between a given pair of nodes. The paths may even overlap along some portion of their length—see Fig. 6.12.

The *diameter* of a network is the length of the “longest shortest path.” That is, among all shortest paths between every pair of nodes in the network for which a path actually exists, the diameter is the length of the longest one.¹⁹ The diameter of the network in Fig. 6.12, for example, is three. The diameter will play a role, for instance, in our proof of the small-world effect for the random graph model in Section 11.7: we will show that the diameter of the network is small in a certain sense, from which it follows that the shortest distance between every pair of nodes is also small (provided the nodes are connected at all).



Any self-intersecting walk must necessarily contain at least one loop (top), and hence can be shortened by removing the loop (bottom).

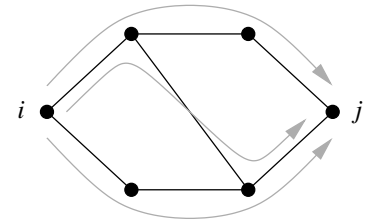


Figure 6.12: There are three shortest paths between nodes i and j in this network, each of length three.

6.12 COMPONENTS

A network need not consist of just a single connected set of nodes. Many networks have two or more separate parts that are disconnected from one

¹⁹If the diameter were merely the length of the longest shortest path then it would be formally infinite in a network with more than one component if we adopted the convention above that nodes connected by no path have infinite distance. One can, however, talk about the diameters of the individual components separately, this being a perfectly well-defined concept whatever convention we adopt for unconnected nodes.

another. For example, the network shown in Fig. 6.13 is divided into two parts, the one on the left having three nodes, the one on the right having four. Such parts are called *components*. There is by definition no path between any pair of nodes in different components. In Fig. 6.13, for instance, there is no path from the node labeled A to the node labeled B.

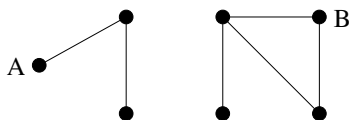


Figure 6.13: A network with two components. There is no path between nodes like A and B that lie in different components.

Technically, a component is a subset of the nodes of a network such that there exists at least one path from each member of that subset to each other member, and such that no other node in the network can be added to the subset while preserving this property. (Subsets like this, to which no other node can be added while preserving a given property, are called *maximal subsets*.) A singleton node that is connected to no others is considered to be a component of size one, and every node belongs to exactly one component. A network in which all nodes belong to the same single component is said to be *connected*. Conversely, a network with more than one component is *disconnected*.

The adjacency matrix of a network with more than one component can be written in block diagonal form, meaning that the non-zero elements of the matrix are confined to square blocks along the diagonal of the matrix, with all other elements being zero:

$$\mathbf{A} = \begin{pmatrix} \boxed{} & 0 & \cdots \\ 0 & \boxed{} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}. \tag{6.26}$$

Note, however, that the node labels must be chosen correctly to give this form. The appearance of blocks in the adjacency matrix relies on the nodes of each component being given sequential labels so that they are grouped together along the axes of the matrix. If the nodes are not grouped in this way the matrix will not be block diagonal and it may be difficult to tell that the network has separate components. There do, however, exist computer algorithms, such as the breadth-first search algorithm described in Section 8.5, that can take a network with arbitrary node labels and quickly determine its components.

6.12.1 COMPONENTS IN DIRECTED NETWORKS

For directed networks the definition of components is more complicated. The situation is worth looking at in some detail, because it assumes some practi-

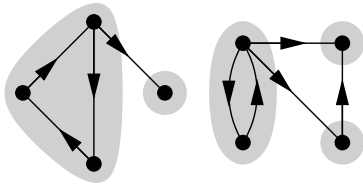


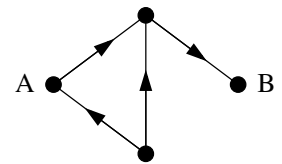
Figure 6.14: Components in a directed network. This network has two weakly connected components of four nodes each, and five strongly connected components (shaded).

cal importance in networks like the World Wide Web. Consider the directed network shown in Fig. 6.14. If we ignore the directed nature of the edges, considering them instead to be undirected, then the network has two components of four nodes each. In the jargon of graph theory these are called *weakly connected components*. Two nodes are in the same weakly connected component if they are connected by one or more paths through the network, where paths are allowed to go either way along any edge.

In many practical situations, however, this is not what we care about. For example, the edges in the World Wide Web are directed hyperlinks that allow web users to surf from one page to another, but only in one direction. This means it is possible to reach one web page from another only if there is a directed path between them, i.e., a path in which we follow edges only in the forward direction. It would be useful to define components for directed networks based on such directed paths, but this raises some problems. It is possible for there to be a directed path from node A to node B but no path back from B to A. Should we then consider A and B to be connected? Are they in the same component or not?

There are various answers one could give to these questions. One possibility is that we define A and B to be connected if and only if there exists a directed path both from A to B and from B to A. In that case, A and B are said to be *strongly connected*. We can define components for a directed network using this definition of connection and these are called *strongly connected components*. Technically, a strongly connected component is a maximal subset of nodes such that there is a directed path in both directions between every pair in the subset. The strongly connected components in the network of Fig. 6.14 are indicated by the shaded regions.

Strongly connected components can consist of just a single node (there are three such components in Fig. 6.14) and every node belongs to exactly one strongly connected component. Note also that every strongly connected component with more than one node must contain at least one cycle. Indeed every node in such a component must belong to at least one cycle, since there is by definition a directed path from that node to every other in the component and a directed path back again, and the two paths together make a cycle. (A



There is a directed path from A to B in this network, but none from B to A.

corollary of this observation is that directed acyclic graphs have no strongly connected components with more than one node, since if they did they wouldn't be acyclic.)

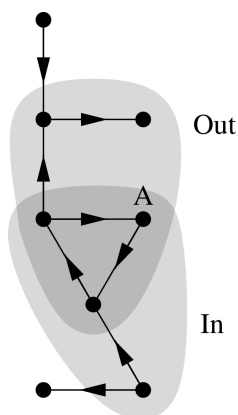


Figure 6.15: In- and out-components. The two shaded regions denote the in- and out-components of node A in this small directed network. The overlap between the two regions is A's strongly connected component.

Strongly and weakly connected components are not the only useful definitions of components in a directed network. On the Web it could be useful to know what pages you can reach by surfing from a given starting point, but you might not care so much whether it's possible to surf back the other way. Considerations of this kind lead us to *out-components*: an out-component is the set of nodes that are reachable via directed paths starting from a specified node A, and including A itself—see Fig. 6.15.

The members of an out-component depend on the choice of the starting node. Choose a different starting node and the set of reachable nodes may change. Thus an out-component is a property of both the network structure and the starting node, and not (as with strongly and weakly connected components) of the network structure alone. This means, among other things, that a node can belong to more than one different out-component. In Fig. 6.16, for instance, we show the out-components of two different starting nodes, A and B. Nodes X and Y belong to both.

A few other points are worth noticing. First, it is self-evident that all the members of the strongly connected component to which node A belongs are also members of A's out-component. Furthermore, any node that is reachable from A is necessarily also reachable from all the other nodes in the strongly connected component. Thus it follows that the out-components of all members of a strongly connected component are identical. It would be reasonable, therefore, to say that out-components really “belong” not to individual nodes, but to strongly connected components.

Note also that while an out-component can have edges to other nodes—nodes not in the out-component—such edges only ever point inward towards the component and never outward (see Fig. 6.16 again for examples). If they pointed outward then the nodes they connected to would by definition be members of the out-component.

Analogous ideas apply also to the nodes *from which* a particular node can be reached. The *in-component* of a specified node A is the set of all nodes from which there is a directed path to A, including A itself (see Fig. 6.15). In-components depend on the choice of the specified node and a node can belong to more than one in-component. But all nodes in the same strongly connected component have the same in-component, and the strongly connected component to which a node belongs is a subset of its in-component. Indeed any

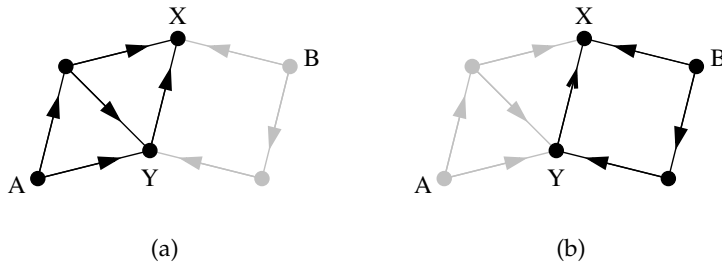


Figure 6.16: Out-components in a directed network. (a) The out-component of node A, which is the set of nodes reachable by directed paths from A. (b) The out-component of node B. Nodes X and Y belong to both out-components.

node that is in both the in-component and the out-component of A is necessarily also in its strongly connected component (since paths exist in both directions) and hence A's strongly connected component is equal to the intersection of its in- and out-components (see Fig. 6.15 again).

6.13 INDEPENDENT PATHS, CONNECTIVITY, AND CUT SETS

There are typically many different ways to walk from one node to another in a network. Even if we restrict ourselves to paths, i.e., self-avoiding walks that never visit the same node twice, there may still be many paths of many different lengths. These paths will usually not be independent however. That is, they will share some nodes or edges, as in Fig. 6.12 for instance. If we restrict ourselves to independent paths, then the number of paths between a given pair of nodes is usually much smaller. Independent paths play an important role in the theory of networks, as we describe in this section.

There are two species of independent path: *edge-independent* and *node-independent*. Two paths connecting a given pair of nodes are edge-independent if they share no edges. Two paths are node-independent if they share no nodes, other than their starting and ending nodes. If two paths are node-independent then they are also edge-independent, but the reverse is not true: it is possible to be edge-independent but not node-independent. For instance, the network shown in Fig. 6.17a has two edge-independent paths from A to B, as denoted by the arrows, but only one node-independent path. The two edge-independent paths are not also node-independent because they share the intermediate node C.

Independent paths are also sometimes called *disjoint paths*, primarily in the mathematical literature. One also sees the terms *edge-disjoint* and *node-disjoint*, describing edge and node independence.

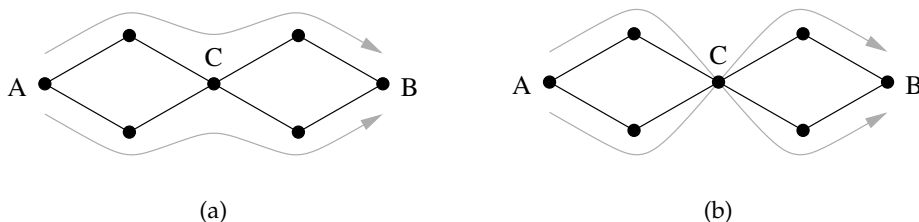


Figure 6.17: Edge independent paths. (a) There are two edge-independent paths from A to B in this figure, as denoted by the arrows, but there is only one node-independent path, because all paths must pass through the center node C. (b) The edge-independent paths are not unique; there are two different ways of choosing the paths from A to B in this case.

The edge- or node-independent paths between two nodes are not necessarily unique. There may be more than one way of choosing a set of independent paths. For instance, Fig. 6.17b shows the same network as Fig. 6.17a, but with the two edge-independent paths chosen a different way, so that they cross over as they pass through the central node C.

It takes only a moment's reflection to convince oneself that there can be only a finite number of independent paths between any two nodes in a network. The number of independent paths (either edge- or node-independent) from A to B cannot exceed A's degree, since every path must leave node A along a different edge. Similarly, the number of paths cannot exceed B's degree either. So the smaller of the degrees of the two nodes gives an upper bound on the number of independent paths. In Fig. 6.17, for instance, there cannot be more than two edge- or node-independent paths between A and B, since both nodes have degree two.

The number of independent paths between a pair of nodes is called the *connectivity* of the nodes.²⁰ If we wish to be explicit about whether we are considering edge- or node-independence, we can refer to *edge* or *node connectivity*. The nodes A and B in Fig. 6.17 have edge connectivity 2 but node connectivity 1 (since there are two edge-independent paths but only one node-independent path).

The connectivity of a pair of nodes can be thought of as a measure of how strongly connected those nodes are. A pair that have only a single independent path between them are, arguably, more tenuously connected than a pair that

²⁰The word "connectivity" is occasionally also used in the networks literature as a synonym for degree, but in the interest of clarity we avoid that usage in this book.

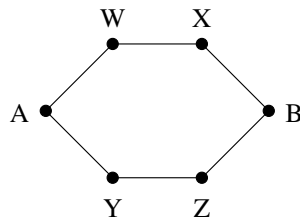
have many paths. This idea is sometimes exploited in the analysis of networks, for instance in algorithmic methods for discovering clusters or communities of strongly linked nodes [181].

Connectivity can also be thought of in terms of “bottlenecks” between nodes. Nodes A and B in Fig. 6.17, for instance, are connected by only one node-independent path because node C forms a bottleneck through which only one path can go. This idea of bottlenecks is formalized by the notion of cut sets as follows.

Consider an undirected network. (In fact the developments here apply equally to directed ones, but for simplicity let us stick with the undirected case for now.) A *cut set*, or more properly a *node cut set*, is a set of nodes whose removal (along with the adjacent edges) will disconnect a specified pair of nodes. For example, the central node C in Fig. 6.17 forms a cut set of size 1 for the nodes A and B. There are also other cut sets for A and B in this network, although all the others are larger than size 1.

An *edge cut set* is the equivalent construct for edges—it is a set of edges whose removal will disconnect a specified pair of nodes.

A *minimum cut set* is the smallest cut set that will disconnect a specified pair of nodes. In Fig. 6.17 the single node C is a minimum node cut set for nodes A and B. A minimum cut set need not be unique. For instance, there are a variety of minimum node cut sets (all of size 2) between the nodes A and B in this network:



$\{W,Y\}$, $\{W,Z\}$, $\{X,Y\}$, and $\{X,Z\}$ are all minimum cut sets for this network. (There are also many different minimum edge cut sets.)

Cut sets were the focus of an important early result in graph theory. *Menger's theorem* says that the size of the minimum cut set between any pair of nodes in a network is equal to the number of independent paths between the same nodes. In other words, the connectivity of a pair of nodes and the number of bottlenecks between them are the same. This theorem applies to both node and edge cut sets and will play an important role when we come to study computer algorithms for analyzing networks, because it allows us to compute the size of a cut set by instead counting independent paths, the latter being a simpler

Algorithms for finding independent paths are discussed in Section 8.7.

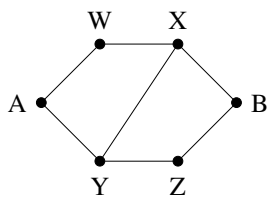


Figure 6.18: Independent paths and cut sets. There are two node-independent paths connecting A and B in this figure, but removing nodes W and Z, one from each path, does not break the connection. Removing X and Y on the other hand, does the job.

operation than the former, as we will see.

Menger's theorem might, at first sight, appear trivial. If there are, say, two node-independent paths between a pair of nodes, then surely we just have to remove one node from each path to sever the connection? Upon further reflection, however, we can see that this doesn't always work. If we remove the wrong nodes from the paths then we don't disconnect the ends—see Fig. 6.18 for an example. Menger's theorem tells us that there always exists *some* set of nodes whose removal will do the job, but they must be the right nodes. Rigorous proof of Menger's theorem is, in fact, not trivial. It was first proved by Karl Menger [330] for the node case, although many other proofs have been given since. A relatively simple one can be found in Ref. [467].

The edge version of Menger's theorem has a further corollary concerning the idea of *maximum flow*. Imagine a set of water pipes in the shape of some network of interest. The edges of the network correspond to the pipes and the nodes to junctions between pipes. And suppose moreover that there is a maximum rate, in terms of volume per unit time, at which water can flow through any pipe—the same maximum for every pipe. What then is the maximum rate of flow from node A to node B through the network as a whole? The answer is that the maximum flow is equal to the number of edge-independent paths times the maximum flow along a single pipe.

This result is called the *max-flow/min-cut theorem*, for the special case in which each pipe can carry the same fixed flow. (There is a more general form that applies when the pipes have different capacities, which we look at in the following section.) The theorem follows straightforwardly from Menger's theorem. The maximum flow must be at least as great as the number of edge-independent paths, since we can simply send one unit of flow along each path. But at the same time it can be no greater than the size of the minimum edge cut set, since if we remove the edges in the cut set we cut off all flow, and each edge in the cut set carries one unit of flow. Since the number of independent paths and the size of the cut set are equal, the result then follows.

Thus, in combination, Menger's theorem and the max-flow/min-cut theorem tell us that for a pair of nodes in an undirected network three quantities are all numerically equal to each other: the edge connectivity of the pair (i.e., the number of edge-independent paths connecting them), the size of the minimum edge cut set (i.e., the number of edges that must be removed to disconnect them), and the maximum flow between the nodes expressed as a multiple of the maximum flow along each individual edge. Although we have stated these results for the undirected case, nothing in any of the proofs demands

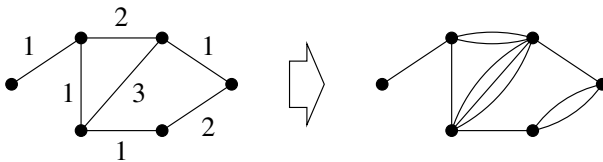
an undirected network, and these three quantities are also equal for directed networks.

6.13.1 MAXIMUM FLOWS AND CUT SETS ON WEIGHTED NETWORKS

As discussed in Section 6.3, networks can have weights on their edges that indicate that some edges are stronger or more prominent than others. In some cases these weights represent capacities of the edges to carry a flow of some kind. For example, they might represent traffic capacity of the roads in a road network or data capacity of Internet lines. We can ask questions about network flows on such networks similar to those we asked in the previous section, but with the added twist that different edges can now have different capacities. For example, we can ask what the maximum possible flow is between a specified pair of nodes. We can also ask about cut sets. An edge cut set for a weighted network is defined just as in the unweighted case to be a set of edges whose removal would disconnect the specified pair of nodes. A *minimum* edge cut set is defined as being a cut set such that the sum of the weights on the edges of the set has the minimum possible value.

Maximum flows and minimum cut sets on weighted networks are related by the general form of the max-flow/min-cut theorem, which says that the maximum flow between a given pair of nodes in a network is equal to the sum of the weights on the edges of the minimum edge cut set that separates the same pair of nodes.

One way to see why this is true is to look at it in terms of the equivalence between weighted networks and multigraphs mentioned in Section 6.3. Consider the special case in which the capacities of all the edges in our network are integers. We can then transform our network by replacing each edge of integer capacity k by k parallel edges of capacity 1, like this:



It is clear that the maximum flow between any two nodes in the transformed network is the same as that between the corresponding nodes in the original. At the same time the transformed network now has the form of an unweighted network of the type considered in Section 6.13, and hence the maximum flow in the original network is equal to the size of the minimum edge cut set in the transformed network.

We further note that the minimum cut set in the transformed network must include either all or none of the parallel edges between any adjacent pair of nodes; there is no point cutting one such edge unless you cut all the others as well. With this constraint, there is a one-to-one correspondence between cut sets on the original network and the transformed network, with corresponding cut sets necessarily having the same total weight. Hence the minimum cut set on the weighted network has the same weight as the minimum cut set on the transformed network and so the minimum cut and maximum flow are equal on the original network.

This demonstrates the theorem for the case of integer edge weights. It can be extended to the non-integer case simply by making the units in which we measure the weights smaller. In the limit where the units become arbitrarily small, any weight can be represented as an integer number of units and the argument above can be applied. Hence the max-flow/min-cut theorem must be generally true for any set of weights.

For an alternate, first-principles proof not using integer edge weights see, for instance, Ahuja *et al.* [9].

There exist efficient computer algorithms for calculating maximum flows on weighted networks, so the max-flow/min-cut theorem allows us to calculate minimum cuts efficiently also, and this is now the standard way of performing such calculations.²¹

6.14 THE GRAPH LAPLACIAN

Section 6.2 introduced the adjacency matrix, which captures the entire structure of a network and whose matrix properties can tell us a variety of useful things. The adjacency matrix, however, is not the only matrix representation of a network. There are several others, including the modularity matrix, the non-backtracking matrix, and the graph Laplacian. Of these, the graph Laplacian is certainly the best known and most widely used.

The graph Laplacian for a simple undirected, unweighted network is an $n \times n$ symmetric matrix \mathbf{L} with elements

$$L_{ij} = \begin{cases} k_i & \text{if } i = j, \\ -1 & \text{if } i \neq j \text{ and there is an edge between nodes } i \text{ and } j, \\ 0 & \text{otherwise,} \end{cases} \quad (6.27)$$

where k_i is the degree of node i , as previously. Another way to write the same

²¹Another (slightly surprising) computational use of the max-flow/min-cut theorem is for finding ground states of the random-field Ising model [373], an interesting case of cross-fertilization from network theory to physics. It is relatively common for physics ideas to find application in network theory, but the reverse is rarer.

thing would be

$$L_{ij} = k_i \delta_{ij} - A_{ij}, \quad (6.28)$$

where A_{ij} is an element of the adjacency matrix and δ_{ij} is the Kronecker delta, which is 1 if $i = j$ and 0 otherwise. Alternatively, we can write \mathbf{L} in matrix form as

$$\mathbf{L} = \mathbf{D} - \mathbf{A}, \quad (6.29)$$

where \mathbf{D} is the diagonal matrix with the node degrees along its diagonal:

$$\mathbf{D} = \begin{pmatrix} k_1 & 0 & 0 & \cdots \\ 0 & k_2 & 0 & \cdots \\ 0 & 0 & k_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}. \quad (6.30)$$

All of these are equivalent definitions of the graph Laplacian.

One can also write a graph Laplacian for weighted networks: one simply replaces the adjacency matrix with the weighted adjacency matrix of Section 6.3 (which has the weights in the matrix elements) and the degree k_i of a node by the sum $\sum_j A_{ij}$ of the relevant matrix elements. One can also treat multigraphs in the same way. There is, however, no natural extension of the graph Laplacian to networks with self-edges or, more importantly, to directed networks. The Laplacian is only useful for the undirected case.

The graph Laplacian crops up in a surprisingly diverse set of situations, including in the theory of random walks on networks, dynamical systems, diffusion, resistor networks, graph visualization, and graph partitioning. In the following sections we look briefly at some of these applications.

6.14.1 GRAPH PARTITIONING

Graph partitioning is the task of dividing the nodes of a network into a set of groups of given sizes so as to minimize the number of edges running between the groups. It arises, for instance, in parallel computing, where you want to divide up a calculation into smaller sub-calculations that can be assigned to several different computers or CPUs, while minimizing the amount of data that will have to be sent back and forth between the CPUs (since transmitting data is usually a relatively cumbersome process that can slow down the whole computation).

Consider the simplest version of graph partitioning, the division of the nodes of a network into just two groups, which we will call group 1 and group 2. The number of edges R running between the two groups, also called

the *cut size*, is given by

$$R = \frac{1}{2} \sum_{\substack{i, j \text{ in} \\ \text{different} \\ \text{groups}}} A_{ij}, \quad (6.31)$$

where the factor of $\frac{1}{2}$ compensates for the fact that every pair of nodes is counted twice in the sum. (For instance, we count nodes 1 and 2 separately from nodes 2 and 1.)

We define a set of quantities s_i , one for each node i , which represent the division of the network thus:²²

$$s_i = \begin{cases} +1 & \text{if node } i \text{ belongs to group 1,} \\ -1 & \text{if node } i \text{ belongs to group 2.} \end{cases} \quad (6.32)$$

Then

$$\frac{1}{2}(1 - s_i s_j) = \begin{cases} 1 & \text{if } i \text{ and } j \text{ are in different groups,} \\ 0 & \text{if } i \text{ and } j \text{ are in the same group,} \end{cases} \quad (6.33)$$

which allows us to rewrite Eq. (6.31) as

$$R = \frac{1}{4} \sum_{ij} A_{ij}(1 - s_i s_j), \quad (6.34)$$

with the sum now over all values of i and j . The first term in the sum is

$$\sum_{ij} A_{ij} = \sum_i k_i = \sum_i k_i s_i^2 = \sum_{ij} k_i \delta_{ij} s_i s_j, \quad (6.35)$$

where k_i is the degree of node i as previously, δ_{ij} is the Kronecker delta, and we have made use of the fact that $\sum_j A_{ij} = k_i$ (see Eq. (6.12)) and $s_i^2 = 1$ (since $s_i = \pm 1$). Substituting back into Eq. (6.34) we then find that

$$R = \frac{1}{4} \sum_{ij} (k_i \delta_{ij} - A_{ij}) s_i s_j = \frac{1}{4} \sum_{ij} L_{ij} s_i s_j, \quad (6.36)$$

where $L_{ij} = k_i \delta_{ij} - A_{ij}$ is the ij th element of the graph Laplacian matrix—see Eq. (6.28).

Equation (6.36) can be written in matrix form as

$$R = \frac{1}{4} \mathbf{s}^T \mathbf{L} \mathbf{s}, \quad (6.37)$$

²²A physicist would call the variables s_i “Ising spins,” and indeed the graph partitioning problem is equivalent to finding the ground state of a certain type of Ising model in which the spins live on the nodes of the network.

where \mathbf{s} is the vector with elements s_i . This expression gives us a matrix formulation of the graph partitioning problem. The matrix \mathbf{L} specifies the structure of our network, the vector \mathbf{s} defines a division of that network into groups, and our goal is to find the vector \mathbf{s} that minimizes the cut size (6.37) for given \mathbf{L} . This matrix formulation leads directly to one of the standard computational methods for solving the graph partitioning problem, *spectral partitioning*, which makes use of the eigenvectors of the graph Laplacian to rapidly find good divisions of the network [177,391].

6.14.2 NETWORK VISUALIZATION

We have seen many pictures of networks in this book. Some of them, like the picture of the Internet on page 2 or the picture of a food web on page 6, depict large and complicated networks that would be difficult to make sense of if the pictures were not carefully laid out to make the network structure as clear as possible. The generation of network visualizations like these is the domain of specialized software packages, whose workings are outside the scope of this book. However, it is interesting to ask, broadly, what is it that characterizes a good visualization of a network?

One answer is that a good visualization is one where the lengths of most edges in the network, as drawn on the page, are short. Consider, for instance, Fig. 6.19, which shows two different pictures of the same network. In Fig. 6.19a the nodes are placed at random on the page, which means that some edges are short but many are relatively long—there are many edges that run clear across the picture from one side to the other. The net result is that the edges are a mess, crossing over one another, getting in each other’s way, and generally making it hard to see which nodes are connected to which. In Fig. 6.19b, on the other hand, the network is laid out so that connected pairs of nodes are (by and large) placed close together and the lengths of the edges are short. This results in a much clearer picture that makes the network structure easier to see.

Suppose then that we have an undirected, unweighted network that we want to lay out on the page. Real network images are two-dimensional but for the sake of simplicity let us consider a one-dimensional case for now, so that the position of node i in our layout is a simple scalar x_i . Our goal is to choose the positions so as to minimize the lengths of edges, which we could do in various ways, but the standard approach is to minimize the sum of the squares of the lengths as follows.

The distance between nodes i and j in our simple one-dimensional model is $|x_i - x_j|$ and the squared distance is $(x_i - x_j)^2$. The sum Δ^2 of the squared

Software packages for network visualization and analysis are discussed in more detail in Chapter 8.

The problem of creating a good visualization of a network is closely related to the theory of graph embeddings and latent spaces, which we examine in Section 14.7.4.

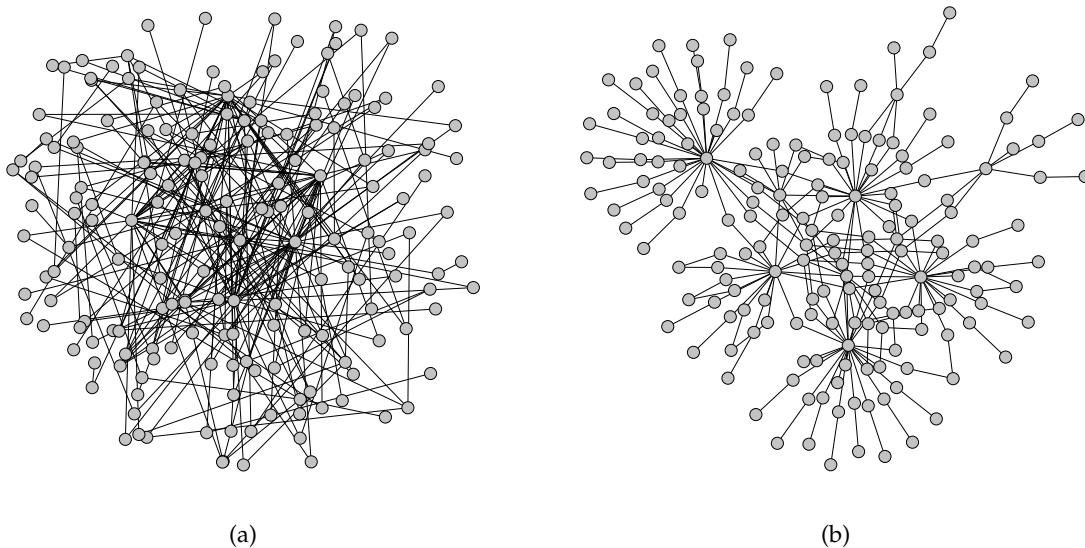


Figure 6.19: Two visualizations of the same network. In (a) nodes are placed randomly on the page, while in (b) nodes are placed using a network layout algorithm that tries to put connected nodes close to one another, meaning that most edges are short.

distances for all node pairs connected by an edge is then

$$\Delta^2 = \frac{1}{2} \sum_{ij} A_{ij} (x_i - x_j)^2, \quad (6.38)$$

where the matrix element A_{ij} ensures that only connected pairs are counted, and the extra factor of $\frac{1}{2}$ compensates for the fact that every pair of nodes appears twice in the sum.

Expanding this expression, we have

$$\begin{aligned} \Delta^2 &= \frac{1}{2} \sum_{ij} A_{ij} (x_i^2 - 2x_i x_j + x_j^2) = \frac{1}{2} \left[\sum_i k_i x_i^2 - 2 \sum_{ij} A_{ij} x_i x_j + \sum_j k_j x_j^2 \right] \\ &= \sum_{ij} (k_i \delta_{ij} - A_{ij}) x_i x_j = \sum_{ij} L_{ij} x_i x_j, \end{aligned} \quad (6.39)$$

where L_{ij} is an element of the graph Laplacian again, and we have made use of the fact that $\sum_j A_{ij} = k_i$ in the second equality (see Eq. (6.12)).

Equation (6.39) can be written in matrix notation as

$$\Delta^2 = \mathbf{x}^T \mathbf{L} \mathbf{x}, \quad (6.40)$$

where \mathbf{x} is the vector with elements x_i . This expression is similar to Eq. (6.37) and, like that equation, it can form the basis for new computer algorithms, in this case algorithms for generating clear visualizations of networks using the eigenvectors of the graph Laplacian [274]. It also tells us that not all networks can be visualized equally clearly. Starting from Eq. (6.40) we can derive a lower bound on the mean-square length of an edge and hence show that in order for a network to have a good visualization where most edges are short it must have low “algebraic connectivity,” meaning that the gap between the smallest and second smallest eigenvalues of the graph Laplacian must be small—see Section 6.14.5. Thus, merely by inspecting the properties of the Laplacian for a particular network we can say whether it will even be possible to make a good visualization. For some networks, no matter how hard we try, we will never be able to make a clear picture because there is no layout in which the average length of edges is small.

6.14.3 RANDOM WALKS

Another context in which the graph Laplacian arises is in the study of random walks on networks. A *random walk* is a walk across a network created by taking repeated random steps. Starting at any initial node, we choose uniformly at random among the edges attached to that node, move along the chosen edge to the node at its other end, and repeat the process. Random walks are allowed to visit the same node more than once, go along the same edge more than once, or backtrack along an edge just traversed. (*Self-avoiding random walks*, which do none of these things, are also studied sometimes, but we will not discuss them here.) Random walks arise, for instance, in the random-walk sampling method for social networks discussed in Section 4.7 and in the random-walk betweenness measure of Section 7.1.7.

Consider a random walk that starts at a specified node and takes t steps. Let $p_i(t)$ be the probability that the walk is at node i at time t . If the walk is at node j at time $t - 1$, the probability of taking a step along any particular one of the k_j edges attached to j is $1/k_j$, so on an undirected network the probability of being at node i on the next step is given by

$$p_i(t) = \sum_j \frac{A_{ij}}{k_j} p_j(t-1), \quad (6.41)$$

or $\mathbf{p}(t) = \mathbf{A}\mathbf{D}^{-1}\mathbf{p}(t-1)$ in matrix form, where \mathbf{p} is the vector with elements p_i and, as before, \mathbf{D} is the diagonal matrix with the degrees of the nodes down its diagonal, as defined in Eq. (6.30).

In the limit of long time the probability distribution over nodes is given by (6.41) with t set to infinity: $p_i(\infty) = \sum_j A_{ij} p_j(\infty) / k_j$, or in matrix form:

$$\mathbf{p} = \mathbf{A}\mathbf{D}^{-1}\mathbf{p}, \quad (6.42)$$

where \mathbf{p} is shorthand for $\mathbf{p}(\infty)$. Rearranging, this can also be written as

$$(\mathbf{I} - \mathbf{A}\mathbf{D}^{-1})\mathbf{p} = (\mathbf{D} - \mathbf{A})\mathbf{D}^{-1}\mathbf{p} = \mathbf{L}\mathbf{D}^{-1}\mathbf{p} = \mathbf{0}. \quad (6.43)$$

Thus $\mathbf{D}^{-1}\mathbf{p}$ is (any multiple of) an eigenvector of the Laplacian with eigenvalue 0.

On a connected network—one with only a single component—we will see in Section 6.14.5 that there is only one eigenvector of the Laplacian that has eigenvalue zero, the vector $\mathbf{1} = (1, 1, 1, \dots)$ whose elements are all 1. Thus $\mathbf{D}^{-1}\mathbf{p} = a\mathbf{1}$, where a is a constant, or equivalently $\mathbf{p} = a\mathbf{D}\mathbf{1}$, so that $p_i = ak_i$. Thus, on a connected network the probability that a random walk will be found at node i in the limit of long time is simply proportional to the degree of that node. If we choose the value of a so that the probabilities p_i sum to one, we get

$$p_i = \frac{k_i}{\sum_j k_j} = \frac{k_i}{2m}, \quad (6.44)$$

where we have made use of Eq. (6.13).

We employed this result previously in Section 4.7 in our analysis of the random-walk sampling method for social networks. The basic insight behind the result is that nodes with high degree are more likely to be visited by a random walk simply because there are more ways of reaching them.

A further corollary is that in the limit of long time the probability $P(i \rightarrow j)$ of walking along an edge from i to j on any particular step of a random walk is equal to the probability p_i of being at node i in the first place times the probability $1/k_i$ of walking along that particular edge:

$$P(i \rightarrow j) = \frac{k_i}{2m} \times \frac{1}{k_i} = \frac{1}{2m}. \quad (6.45)$$

In other words, on any given step a random walk is equally likely to traverse every edge.

6.14.4 RESISTOR NETWORKS

As a further example of the application of the graph Laplacian, consider a network of resistors, one of the simplest examples of an electrical network. Suppose we have a network in which the edges are identical resistors of resistance R and the nodes are junctions between resistors, as shown in Fig. 6.20,

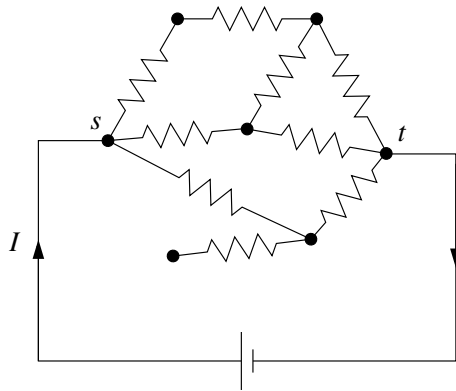


Figure 6.20: A resistor network with an applied voltage. In this network the edges are resistors and the nodes are electrical junctions between them. A voltage is applied between nodes s and t , generating a total current I .

and suppose we apply a voltage between two nodes s and t such that a total current I flows from s to t through the network.

One basic question we could ask about such a network is what the voltage is at any given node. The current flow in the network obeys Kirchhoff's current law, which is essentially a statement that electricity is conserved, so that the net current flowing in or out of any node is zero. Let V_i be the voltage at node i , measured relative to any convenient reference potential. Then Kirchhoff's law says that

$$\sum_j A_{ij} \frac{V_i - V_j}{R} - I_i = 0, \quad (6.46)$$

where I_i represents any current injected into node i by an external current source. In our case this external current is non-zero only for the two nodes s and t connected to the external voltage:

$$I_i = \begin{cases} +I & \text{for } i = s, \\ -I & \text{for } i = t, \\ 0 & \text{otherwise.} \end{cases} \quad (6.47)$$

(In theory there's no reason why one could not impose more complex current source arrangements by applying additional voltages to the network and making more elements I_i non-zero, but let us stick to our simple case in this discussion.)

Noting that $\sum_j A_{ij} = k_i$, Eq. (6.46) can also be written as $k_i V_i - \sum_j A_{ij} V_j = R I_i$ or

$$\sum_j (\delta_{ij} k_i - A_{ij}) V_j = R I_i, \quad (6.48)$$

which in matrix form is

$$\mathbf{L}\mathbf{V} = \mathbf{R}\mathbf{I}, \quad (6.49)$$

where \mathbf{L} is once again the graph Laplacian. This equation is a kind of matrix version of the standard Ohm's law $V = RI$ for a single resistor, and by solving it for \mathbf{V} we can calculate the voltages at every node in the network.

Calculating the behavior of a resistor network might seem like a problem of rather narrow interest, but in fact the connection between the Laplacian and resistor networks has an important and perhaps surprising practical application. It is the basis for the most widely used technique for *graph sparsification*, in which one aims to remove edges from a network while keeping other properties of the network the same [49,435]. Graph sparsification forms the foundation for a range of modern numerical methods for solving large systems of simultaneous linear equations. By representing the equations in the form of a resistor network (in effect the reverse of the operations above, where we took a resistor network and represented it by a set of equations), then sparsifying that network while keeping its electrical properties the same, we can enormously reduce the complexity of the problem, allowing us to solve in seconds systems of equations that previously might have taken hours. Graph sparsification and its application in equation solving is just one example of the many important technological uses of network theory.

6.14.5 PROPERTIES OF THE GRAPH LAPLACIAN

The graph Laplacian has a number of specific properties that are important in many calculations. For instance, it has the property that every row of the matrix sums to zero:

$$\sum_j L_{ij} = \sum_j (k_i \delta_{ij} - A_{ij}) = k_i - k_i = 0, \quad (6.50)$$

where we have made use of the fact that $\sum_j A_{ij} = k_i$ —see Eq. (6.12). Similarly every column of the matrix also sums to zero.

Of particular interest are the eigenvalues of the graph Laplacian. Since the Laplacian is a real symmetric matrix, it necessarily has real eigenvalues. But we can say more than this: all the eigenvalues of the Laplacian are also non-negative.

Let λ be any eigenvalue of the graph Laplacian and let \mathbf{v} be the corresponding eigenvector, unit normalized so that $\mathbf{v}^T \mathbf{v} = 1$. Then $\mathbf{L}\mathbf{v} = \lambda \mathbf{v}$ and

$$\mathbf{v}^T \mathbf{L}\mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda. \quad (6.51)$$

Following the same line of argument we used in Eqs. (6.38) to (6.40) we can write

$$\begin{aligned} \sum_{ij} A_{ij}(v_i - v_j)^2 &= \sum_{ij} A_{ij}(v_i^2 - 2v_i v_j + v_j^2) \\ &= \sum_i k_i v_i^2 - 2 \sum_{ij} A_{ij} v_i v_j + \sum_j k_j v_j^2 \\ &= 2 \sum_{ij} (k_i \delta_{ij} - A_{ij}) v_i v_j = 2 \sum_{ij} L_{ij} v_i v_j = 2 \mathbf{v}^T \mathbf{L}\mathbf{v}. \end{aligned} \quad (6.52)$$

And combining (6.51) and (6.52) we then get

$$\lambda = \frac{1}{2} \sum_{ij} A_{ij}(v_i - v_j)^2 \geq 0. \quad (6.53)$$

Thus all eigenvalues of the Laplacian are non-negative.

While the eigenvalues cannot be negative, however, they can be zero, and in fact the Laplacian always has at least one zero eigenvalue. As we have seen, every row of the matrix sums to zero, which means that the vector $\mathbf{1} = (1, 1, 1, \dots)$ is always an eigenvector of the Laplacian with eigenvalue zero: $\mathbf{L}\mathbf{1} = 0$. (It is not a properly normalized eigenvector. The properly normalized vector would be $(1, 1, 1, \dots)/\sqrt{n}$.) Since there are no negative eigenvalues, this is the lowest of the eigenvalues of the Laplacian.

The presence of a zero eigenvalue implies, among other things, that the Laplacian has no inverse: the determinant of a matrix is the product of its eigenvalues, and hence the determinant of the Laplacian is always zero, so the matrix is singular.

The Laplacian can have more than one zero eigenvalue. Consider, for instance, a network that is divided into c different components of sizes n_1, \dots, n_c and let us number the nodes of the network so that the first n_1 nodes are those of the first component, the next n_2 are those of the second component, and so forth. With this choice the Laplacian of the network is block diagonal, looking something like this:

See the discussion of block diagonal matrices in Section 6.12.

$$\mathbf{L} = \begin{pmatrix} \boxed{} & 0 & \cdots \\ 0 & \boxed{} & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (6.54)$$

What is more, each individual block in the matrix is itself the Laplacian of the corresponding component: it has the degrees of the nodes in that component along its diagonal and -1 in each position corresponding to an edge. This implies that each block has its own eigenvector $(1, 1, 1, \dots)$ with eigenvalue zero, which in turn tells us that there must be at least c different (linearly independent) vectors that are eigenvectors of the full Laplacian \mathbf{L} with eigenvalue zero: the vectors that have ones in the positions corresponding to the nodes in a single component and zeros everywhere else. For instance, the vector

$$\mathbf{v} = (\underbrace{1, 1, 1, \dots}_{n_1 \text{ ones}}, \underbrace{0, 0, 0, \dots}_{\text{zeros}}), \quad (6.55)$$

is an eigenvector with eigenvalue zero. Thus, in a network with c components there are always at least c zero eigenvalues.

Conversely, one can also show that if the network has only one component then the graph Laplacian has only a single zero eigenvalue. To see this, consider an eigenvector \mathbf{v} with eigenvalue zero. Equation (6.53) tells us that for this vector $\sum_{ij} A_{ij}(v_i - v_j)^2 = 0$, which can only be true if $v_i = v_j$ at opposite ends of every edge. If the network has only one component, however, so that we can get from any node to any other by walking along a suitable sequence of edges, this implies that v_i must have the same value at every node, in which case \mathbf{v} is just a multiple of the vector $\mathbf{1}$. In other words, in a network with only one component there is only one eigenvector with eigenvalue zero, the vector $\mathbf{1}$ (or multiples of it). All other eigenvectors must have non-zero eigenvalues.

To put this another way, if a network has only one component then the second smallest eigenvalue will be non-zero. At the same time, as we have said, a network with more than one component will have more than one zero eigenvalue, meaning that the second smallest one is zero. Thus, the second smallest eigenvalue is non-zero if and only if the network is connected—if it consists of a single component. The second smallest eigenvalue of the Laplacian is called the *algebraic connectivity* of the network or the *spectral gap*. It plays an important role in a number of areas of network theory.

It is a straightforward extension of the same arguments to show that the number of zero eigenvalues of the Laplacian is in fact always exactly equal to the number of components in the network. As described earlier, the Laplacian for a network with more than one component is block diagonal, with each block taking the form of the Laplacian for the corresponding component. Each such block Laplacian necessarily has only one zero eigenvalue because the component it describes is, by definition, connected. Hence all blocks together contribute exactly as many zero eigenvalues to the Laplacian of the complete network as there are blocks. For a more detailed proof see, for example, West [467].

EXERCISES

6.1 Which word or words from the following list describe each of the five networks below: *directed, undirected, cyclic, acyclic, approximately acyclic, planar, approximately planar, tree, approximate tree*.

- a) The Internet, at the level of autonomous systems
- b) A food web
- c) The stem and branches of a plant
- d) A spider web
- e) A complete clique of four nodes

Give one real-life example of each of the following types of networks, not including the five examples above:

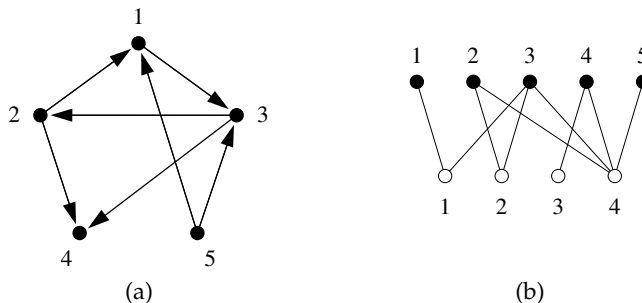
- f) An acyclic (or approximately acyclic) directed network
- g) A cyclic directed network
- h) A tree (or approximate tree)
- i) A planar (or approximately planar) network
- j) A bipartite network

Describe briefly one empirical technique that could be used to measure the structure of each of the following networks (i.e., to fully determine the positions of all the edges):

- k) The World Wide Web
- l) A citation network of scientific papers
- m) A food web
- n) A network of friendships between a group of co-workers
- o) A power grid

6.2 A simple network consists of n nodes in a single component. What is the maximum possible number of edges it could have? What is the minimum possible number of edges it could have? Explain briefly how you arrive at your answers.

6.3 Consider the following two networks:

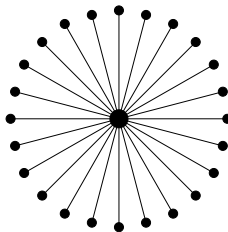


Network (a) is directed. Network (b) is undirected but bipartite. Write down:

- a) The adjacency matrix of network (a);
 - b) The incidence matrix of network (b);
 - c) The projection matrix (Eq. (6.10)) for the projection of network (b) onto its black nodes.
- 6.4 Let A be the adjacency matrix of an undirected network and $\mathbf{1}$ be the column vector whose elements are all 1. In terms of these quantities write expressions for:
- a) The vector \mathbf{k} whose elements are the degrees k_i of the nodes;
 - b) The number m of edges in the network;
 - c) The matrix \mathbf{N} whose element N_{ij} is equal to the number of common neighbors of nodes i and j ;
 - d) The total number of triangles in the network, where a triangle means three nodes, each connected by edges to both of the others.

- 6.5 Demonstrate the following for undirected networks:
- a) A 3-regular graph must have an even number of nodes.
 - b) The average degree of a tree is strictly less than 2.
 - c) Consider any three nodes A, B, and C in a network. The edge connectivity of A and B is x . The edge connectivity of B and C is y , with $y < x$. What is the edge connectivity of A and C, and why?

6.6 A “star graph” consists of a single central node with $n - 1$ other nodes connected to it thus:



What is the largest (most positive) eigenvalue of the adjacency matrix of this network?

6.7 Consider an acyclic directed network of n nodes, labeled $i = 1 \dots n$, and suppose that the labels are assigned in the manner of Fig. 6.3 on page 111, such that all edges run from nodes with higher labels to nodes with lower.

- a) Write down an expression for the total number of ingoing edges at nodes $1 \dots r$ and another for the total number outgoing at nodes $1 \dots r$, in terms of the in- and out-degrees k_i^{in} and k_i^{out} of the nodes.
- b) Hence find an expression for the total number of edges running to nodes $1 \dots r$ from nodes $r + 1 \dots n$.
- c) Hence or otherwise show that in any acyclic network the in- and out-degrees must satisfy

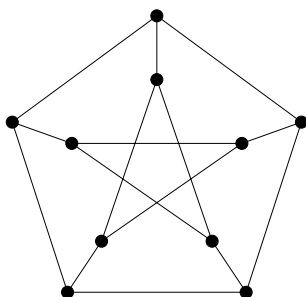
$$k_r^{\text{in}} \leq \sum_{i=r+1}^n (k_i^{\text{out}} - k_i^{\text{in}}), \quad k_{r+1}^{\text{out}} \leq \sum_{i=1}^r (k_i^{\text{in}} - k_i^{\text{out}}),$$

for all r .

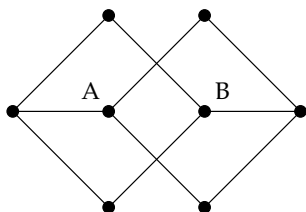
6.8 Consider a bipartite network, with its two types of nodes, and suppose that there are n_1 nodes of type 1 and n_2 nodes of type 2. Show that the mean degrees c_1 and c_2 of the two types are related by

$$c_2 = \frac{n_1}{n_2} c_1.$$

6.9 Using Kuratowski's theorem, prove that this network is not planar:

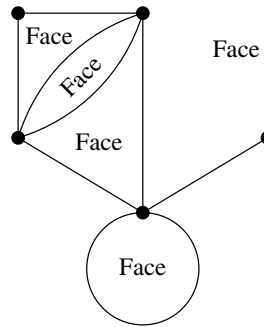


6.10 Give a proof, visual or otherwise, that the edge connectivity of nodes A and B in this network is 2:



Hint: A correct proof must show both that the connectivity is at least 2 and that it is no more than 2.

6.11 Consider a connected planar network with n nodes and m edges. Let f be the number of “faces” of the network, i.e., areas bounded by edges when the network is drawn in planar form. The “outside” of the network, the area extending to infinity on all sides, is also considered a face. The network can have multiedges and self-edges:



- Write down the values of n , m , and f for a network with a single node and no edges.
 - How do n , m , and f change when we add a single node to the network along with a single edge attaching it to another node?
 - How do n , m , and f change when we add a single edge between two extant nodes (or a self-edge attached to just one node), in such a way as to maintain the planarity of the network?
 - Hence by induction prove a general relation between n , m , and f for all connected planar networks.
 - Now suppose that our network is simple (i.e., it contains no multiedges or self-edges). Show that the mean degree c of such a network is strictly less than six.
- 6.12 Consider the set of all paths from node s to node t on an undirected network with adjacency matrix \mathbf{A} . Let us give each path a weight equal to α^r , where α is a constant and r is the length of the path.

- Show that the sum of the weights of all the paths from s to t is given by the st element of the matrix $\mathbf{Z} = (\mathbf{I} - \alpha\mathbf{A})^{-1}$, where \mathbf{I} is the identity matrix.
- What condition must α satisfy for the sum to converge?
- Hence, or otherwise, show that the length ℓ_{st} of the shortest path from s to t , if there is one, is

$$\ell_{st} = \lim_{\alpha \rightarrow 0} \frac{\partial \log Z_{st}}{\partial \log \alpha}.$$

- 6.13 In Section 5.3.1, we gave one possible definition of the trophic level x_i of a species in a (directed) food web as the mean of the trophic levels of the species' prey, plus one.
- Show that x_i , when defined in this way, satisfies

$$x_i = 1 + \frac{1}{k_i^{\text{in}}} \sum_j A_{ij} x_j.$$

- b) This expression does not work for autotrophs—species with no prey—because the corresponding vector element is undefined. Such species are usually given a trophic level of one. Suggest a modification of the calculation that will correctly assign trophic levels to these species, and hence to all species. Thus, show that x_i can be calculated as the i th element of a vector

$$\mathbf{x} = (\mathbf{D} - \mathbf{A})^{-1} \mathbf{D} \cdot \mathbf{1},$$

and specify how the matrix \mathbf{D} is defined.

CHAPTER 7

MEASURES AND METRICS

An introduction to some of the standard measures for quantifying network structure, including measures of centrality, transitivity, and modularity

IF WE HAVE a complete map of a network, of which nodes are connected to which, then in principle we know everything there is to know about its structure. In practice, however, raw network data are not easy for humans to comprehend. For a small network we might be able to make a useful visualization and learn something by inspecting it, but this approach does not work for larger networks, and even for the small ones it can only give us a qualitative feel for the data.

A better approach is to define mathematical measures that capture interesting features of network structure quantitatively, boiling down large volumes of complex structural data into simple numbers that are easy for people to understand. Many such measures have been proposed over the years and in this chapter we look at some of the most widely used. Many of the ideas in this area come from the social sciences, from the discipline of *social network analysis*, which was developed to aid our understanding of social network data such as those described in Chapter 4, and much of the language used to describe these ideas reflects their sociological origin. However, the methods described are now in wide use across many other areas as well, including computer science, physics, statistics, and biology, and they form an important part of the basic network toolbox.

For those interested in traditional social network analysis, introductions can be found in the books by Scott [424] and by Wasserman and Faust [462].

Networks, 2nd edition. Mark Newman, Oxford University Press (2018). © Mark Newman.
DOI: 10.1093/oso/9780198805090.001.0001

7.1 CENTRALITY

A large volume of research on networks has been devoted to the concept of *centrality*. This research addresses the question, “Which are the most important or central nodes in a network?” There are many possible definitions of importance and there are correspondingly many centrality measures for networks. In the following several sections we describe some of the most widely used such measures.

7.1.1 DEGREE CENTRALITY

Perhaps the simplest centrality measure for a node in a network is just its degree, the number of edges connected to it (see Section 6.10). Degree is sometimes called *degree centrality* in the social networks literature, to emphasize its use as a centrality measure. In directed networks, nodes have both an in-degree and an out-degree, and both may be useful as measures of centrality in the appropriate circumstances.

Although degree centrality is a simple centrality measure, it can be very illuminating. In a social network, for instance, it seems reasonable to suppose that individuals who have many friends or acquaintances might have more influence, more access to information, or more prestige than those who have fewer. A non-social network example is the use of citation counts in the evaluation of scientific papers. The number of citations a paper receives from other papers, which is its in-degree in the directed citation network, gives a quantitative measure of how influential the paper has been and is widely used for judging the impact of scientific research.

7.1.2 EIGENVECTOR CENTRALITY

Useful though it is, degree is quite a crude measure of centrality. In effect, it awards a node one “centrality point” for every neighbor it has. But not all neighbors are necessarily equivalent. In many circumstances a node’s importance in a network is increased by having connections to other nodes that *are themselves important*. For instance, you might have only one friend in the world, but if that friend is the president of the United States then you yourself may be an important person. Thus centrality is not only about how many people you know but also who you know.

Eigenvector centrality is an extension of degree centrality that takes this factor into account. Instead of just awarding one point for every network neighbor a node has, eigenvector centrality awards a number of points *proportional to the centrality scores of the neighbors*. This might sound tautological—in order to

work out the score of every node, I need to know the score of every node. But in fact it is straightforward to calculate the scores with just a little work.

Consider an undirected network of n nodes. The eigenvector centrality x_i of node i is defined to be proportional to the sum of the centralities of i 's neighbors, so that

$$x_i = \kappa^{-1} \sum_{\substack{\text{nodes } j \text{ that are} \\ \text{neighbors of } i}} x_j, \quad (7.1)$$

where we have called the constant of proportionality κ^{-1} for reasons that will become clear. For the moment we will leave the value of κ arbitrary—we will choose a value shortly.

With eigenvector centrality defined in this way, a node can achieve high centrality either by having a lot of neighbors with modest centrality, or by having a few neighbors with high centrality (or both). This seems natural: you can be influential either by knowing a lot of people, or by knowing just a few people if those few are themselves influential.

An alternative way to write Eq. (7.1) is to make use of the adjacency matrix:

$$x_i = \kappa^{-1} \sum_{j=1}^n A_{ij} x_j. \quad (7.2)$$

Note that the sum is now over all nodes j —the factor of A_{ij} ensures that only the terms for nodes that are neighbors of i contribute to the sum. This formula can also be written in matrix notation as $\mathbf{x} = \kappa^{-1} \mathbf{A}\mathbf{x}$, or equivalently

$$\mathbf{A}\mathbf{x} = \kappa\mathbf{x}, \quad (7.3)$$

where \mathbf{x} is the vector with elements equal to the centrality scores x_i . In other words, \mathbf{x} is an eigenvector of the adjacency matrix.

This doesn't completely fix the centrality scores, however, since there are n different eigenvectors of the $n \times n$ adjacency matrix. Which eigenvector should we use? Assuming we want our centrality scores to be non-negative, there is only one choice: \mathbf{x} must be the leading eigenvector of the adjacency matrix, i.e., the eigenvector corresponding to the largest (most positive) eigenvalue. We can say this with certainty because of the *Perron–Frobenius theorem*, one of the most famous and fundamental results in linear algebra, which states that for a matrix with all elements non-negative, like the adjacency matrix, there is only one eigenvector that also has all elements non-negative, and that is the leading eigenvector. Every other eigenvector must have at least one negative

element.^{1,2}

So this is the definition of the eigenvector centrality, as first proposed by Bonacich in 1987 [73]: the centrality x_i of node i is the i th element of the leading eigenvector of the adjacency matrix.

This also fixes the value of the constant κ —it must be equal to the largest eigenvalue. The centrality scores themselves are still arbitrary to within a multiplicative constant: if we multiply all elements of \mathbf{x} by any constant, Eq. (7.3) is unaffected. In most cases this doesn't matter much. Usually the purpose of centrality scores is to pick out the most important nodes in a network or to rank nodes from most to least important, so it is only the relative scores of different nodes that matter, not the absolute numbers. If we wish, however, we can normalize the centralities by, for instance, requiring that they sum to n (which ensures that average centrality stays constant as the network gets larger).

We have defined the eigenvector centrality here for the case of an undirected network. In theory it can be calculated for directed networks too, but it works best for the undirected case. In the directed case other complications arise. First of all, a directed network has an adjacency matrix that is, in general, asymmetric (see Section 6.4). This means it has two sets of eigenvectors, the left eigenvectors and the right eigenvectors, and hence two leading eigenvectors. Which of the two should we use to define the centrality? In most cases the

¹Technically, this result is only true for connected networks, i.e., networks with only one component. If a network has more than one component then there is one eigenvector with non-negative elements for each component. This doesn't pose a practical problem though: one can simply split the network into its components and calculate the eigenvector centrality separately for each one, which again guarantees that there is only one vector with all elements non-negative.

²A detailed discussion and proof of the Perron–Frobenius can be found, for example, in the books by Meyer [331] and Strang [440]. The basic intuition behind it is simple though. Suppose we take a random vector $\mathbf{x}(0)$ and multiply it repeatedly by a symmetric matrix \mathbf{A} that has all elements non-negative. After t multiplications we get a vector $\mathbf{x}(t) = \mathbf{A}^t \mathbf{x}(0)$. Now let us write $\mathbf{x}(0)$ as a linear combination $\mathbf{x}(0) = \sum_i c_i \mathbf{v}_i$ of the eigenvectors \mathbf{v}_i of \mathbf{A} , for some appropriate choice of constants c_i . Then

$$\mathbf{x}(t) = \mathbf{A}^t \mathbf{x}(0) = \mathbf{A}^t \sum_i c_i \mathbf{v}_i = \sum_i c_i \kappa_i^t \mathbf{v}_i = \kappa_1^t \sum_i c_i \left(\frac{\kappa_i}{\kappa_1} \right)^t \mathbf{v}_i,$$

where κ_i are the eigenvalues of \mathbf{A} and κ_1 is the eigenvalue of largest magnitude. Since $|\kappa_i/\kappa_1| < 1$ for all $i \neq 1$, all terms in the sum other than the first decay exponentially as t becomes large, and hence in the limit $t \rightarrow \infty$ we get $\mathbf{x}(t)/\kappa_1^t \rightarrow c_1 \mathbf{v}_1$. In other words, the limiting vector is simply proportional to the leading eigenvector of the matrix.

But now suppose we choose our initial vector $\mathbf{x}(0)$ to have only non-negative elements. Since all elements of the adjacency matrix are also non-negative, multiplication by \mathbf{A} can never introduce any negative elements into the vector and $\mathbf{x}(t)$ must have all elements non-negative for all values of t . Thus the leading eigenvector of \mathbf{A} must also have all elements non-negative. As a corollary, this also implies that κ_1 must be positive, since $\mathbf{A}\mathbf{x} = \kappa_1\mathbf{x}$ has no solutions for negative κ_1 if both \mathbf{A} and \mathbf{x} have only non-negative elements.

correct answer is to use the right eigenvector. The reason is that centrality in directed networks is usually bestowed by other nodes that point towards you, rather than by you pointing to others. On the World Wide Web, for instance, it is a good indication of the importance of a web page that it is pointed to by many other important web pages. On the other hand, the fact that a page might itself point to important pages is neither here nor there. Anyone can set up a page that points to a thousand others, but that does not make the page important.³ Similar considerations apply also to citation networks and other directed networks. Thus the correct definition of eigenvector centrality for a node i in a directed network makes it proportional to the centralities of the nodes that point to it:

$$x_i = \kappa^{-1} \sum_j A_{ij} x_j, \quad (7.4)$$

which gives $\mathbf{Ax} = \kappa\mathbf{x}$ in matrix notation, where \mathbf{x} is the right leading eigenvector.

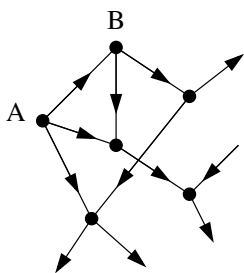


Figure 7.1: A portion of a directed network. Node A in this network has only outgoing edges and hence will have eigenvector centrality zero. Node B has outgoing edges and one ingoing edge, but the ingoing one originates at A, and hence node B will also have centrality zero.

However, there are still problems with this definition. Consider Fig. 7.1. Node A in this figure is connected to the rest of the network, but has only outgoing edges and no ingoing ones. Such a node will always have eigenvector centrality zero because all terms in the sum in Eq. (7.4) are zero. This might not seem to be a problem: perhaps a node that no one points to *should* have centrality zero. But then consider node B. Node B has one ingoing edge, but that edge originates at node A, and hence B also has centrality zero—all terms in the sum in Eq. (7.4) are again zero. Taking this argument further, we see that a node may be pointed to by others that themselves are pointed to by many more, and so on through many generations, but if the trail ends at a node or nodes that have in-degree zero, it is all for nothing—the final value of the centrality will still be zero.

In mathematical terms, only nodes that are in a strongly connected component of two or more nodes, or the out-component of such a strongly connected component, can have non-zero eigenvector centrality.⁴ In many cases, however, it is appropriate for nodes with high in-degree to have high centrality even if they are not in a strongly-connected component or its out-component. Web pages with many links, for instance, can reasonably be considered important even if they are not

³Arguably, this is not entirely true. Web pages that point to many others are often directories of one sort or another and can be useful as starting points for web surfing. This, however, is a different kind of importance from that highlighted by the eigenvector centrality and a different, complementary centrality measure is needed to quantify it. See Section 7.1.5.

⁴For the left eigenvector it would be the in-component.

in a strongly connected component. Recall also that acyclic networks, such as citation networks, have no strongly connected components of more than one node (see Section 6.12.1), so all nodes will have centrality zero, making the eigenvector centrality completely useless for acyclic networks.

There are a number of variants of eigenvector centrality that address these problems. In the next two sections we discuss two of them: Katz centrality and PageRank.

7.1.3 KATZ CENTRALITY

One solution to the issues of the previous section is the following: we simply give each node a small amount of centrality “for free,” regardless of its position in the network or the centrality of its neighbors. In other words, we define

$$x_i = \alpha \sum_j A_{ij}x_j + \beta, \quad (7.5)$$

where α and β are positive constants. The first term is the normal eigenvector centrality term in which the centralities of the nodes pointing to i are summed, and the second term is the “free” part, the constant extra amount that all nodes receive. By adding this second term, we ensure that even nodes with zero in-degree still get centrality β , and once they have non-zero centrality they can pass it along to the other nodes they point to. This means that any node that is pointed to by many others will have a high centrality, even if it is not in a strongly connected component or out-component.

In matrix terms, Eq. (7.5) can be written

$$\mathbf{x} = \alpha \mathbf{A}\mathbf{x} + \beta \mathbf{1}, \quad (7.6)$$

where $\mathbf{1}$ is the uniform vector $(1, 1, 1, \dots)$. Rearranging for \mathbf{x} , we then find that $\mathbf{x} = \beta(\mathbf{I} - \alpha\mathbf{A})^{-1}\mathbf{1}$. As we have said, we normally don’t care about the absolute magnitude of centrality scores, only about the relative scores of different nodes, so the overall multiplier β is unimportant. For convenience we usually set $\beta = 1$, giving

$$\mathbf{x} = (\mathbf{I} - \alpha\mathbf{A})^{-1}\mathbf{1}. \quad (7.7)$$

This centrality measure was first proposed by Katz in 1953 [258] and we will refer to it as the *Katz centrality*.

The definition of the Katz centrality contains the parameter α , which governs the balance between the eigenvector centrality term in Eq. (7.5) and the constant term. If we wish to make use of the Katz centrality we must first choose a value for this constant. In doing so it is important to understand

that α cannot be arbitrarily large. If we let $\alpha \rightarrow 0$, then only the constant term survives in Eq. (7.5) and all nodes have the same centrality β (which we have set to 1). As we increase α from zero the centralities increase and eventually there comes a point at which they diverge. This happens when $(\mathbf{I} - \alpha\mathbf{A})^{-1}$ diverges in Eq. (7.7), i.e., when $\det(\mathbf{I} - \alpha\mathbf{A})$ passes through zero. Rewriting this condition as

$$\det(\alpha^{-1}\mathbf{I} - \mathbf{A}) = 0, \quad (7.8)$$

we see that it is simply the characteristic equation whose roots α^{-1} are equal to the eigenvalues of the adjacency matrix.⁵ As α increases, the determinant first crosses zero when $\alpha^{-1} = \kappa_1$, the largest (most positive) eigenvalue of \mathbf{A} , or alternatively when $\alpha = 1/\kappa_1$. Thus, we should choose a value of α less than this if we wish the expression for the centrality to converge.⁶

Beyond this, however, there is little guidance to be had as to the value that α should take. Most researchers have employed values close to the maximum of $1/\kappa_1$, which places the maximum amount of weight on the eigenvector term and the smallest amount on the constant term. This returns a centrality that is numerically quite close to the ordinary eigenvector centrality, but gives small non-zero values to nodes that are not in strongly connected components of size two or more or their out-components.⁷

The Katz centrality provides a solution to the problems encountered with ordinary eigenvector centrality in directed networks. However, there is no reason in principle why one cannot use Katz centrality in undirected networks as well, and there are times when this might be worthwhile. The idea of adding a constant term to the centrality so that each node gets some weight just by virtue of existing is a natural one. It allows a node that has many neighbors to have high centrality regardless of whether those neighbors themselves have high centrality, and this could be useful in some applications.

⁵The determinant of a matrix is equal to the product of the eigenvalues of the matrix. The matrix $x\mathbf{I} - \mathbf{A}$ has eigenvalues $x - \kappa_i$ where κ_i are the eigenvalues of \mathbf{A} , and hence its determinant is $\det(x\mathbf{I} - \mathbf{A}) = (x - \kappa_1)(x - \kappa_2) \dots (x - \kappa_n)$, which is a degree- n polynomial in x with zeros at $x = \kappa_1, \kappa_2, \dots$. Hence the solutions of $\det(x\mathbf{I} - \mathbf{A}) = 0$ give the eigenvalues of \mathbf{A} .

⁶Formally one recovers finite values again when one moves past $1/\kappa_1$ to higher α , but in practice these values are meaningless. The method returns good results only for $\alpha < 1/\kappa_1$.

⁷In fact, the Katz centrality becomes formally equal to the eigenvector centrality in the limit $\alpha \rightarrow 1/\kappa_1$. Moreover, it is equivalent to degree centrality in the limit $\alpha \rightarrow 0$. So the Katz centrality includes both these other measures as special cases and interpolates between them for intermediate values of α . See Exercise 7.3 for more details.

7.1.4 PAGERANK

The Katz centrality of the previous section has one potentially undesirable feature. If a node with high Katz centrality has edges pointing to many others then all of those others also get high centrality. A high-centrality node pointing to one million others gives all one million of them high centrality. One could argue that this is not always appropriate. In many cases it means less if a node is only one among many that are pointed to. The centrality gained by virtue of receiving an edge from a prestigious node is diluted by being shared with so many others. For instance, websites like *Amazon* or *eBay* link to the web pages of thousands of manufacturers and sellers; if I'm selling something on Amazon it might link to me. Amazon is an important website, and would have high centrality by any sensible measure, but should I therefore be considered important by association? Most people would say not: I am only one of many that Amazon links to and its contribution to the centrality of my page will get diluted as a result.

We can allow for this by defining a variant of the Katz centrality in which the centrality I derive from my network neighbors is proportional to their centrality *divided by their out-degree*. Then nodes that point to many others pass only a small amount of centrality on to each of those others, even if their own centrality is high.

In mathematical terms this centrality is defined by

$$x_i = \alpha \sum_j A_{ij} \frac{x_j}{k_j^{\text{out}}} + \beta. \quad (7.9)$$

This gives problems, however, if there are nodes in the network with out-degree $k_j^{\text{out}} = 0$. For such nodes the corresponding term in the sum in Eq. (7.9) is indeterminate—it is equal to zero divided by zero (because A_{ij} is always zero if j has no outgoing edges). This problem is easily fixed however. It is clear that nodes with no out-going edges should contribute zero to the centrality of any other node, which we can contrive by artificially setting $k_j^{\text{out}} = 1$ for all such nodes. (In fact, we could set k_j^{out} to any non-zero value and the calculation would give the same answer.)

In matrix terms, Eq. (7.9) is then

$$\mathbf{x} = \alpha \mathbf{A} \mathbf{D}^{-1} \mathbf{x} + \beta \mathbf{1}, \quad (7.10)$$

with $\mathbf{1}$ being again the vector $(1, 1, 1, \dots)$ and \mathbf{D} being the diagonal matrix with elements $D_{ii} = \max(k_i^{\text{out}}, 1)$. Rearranging, we find that $\mathbf{x} = \beta (\mathbf{I} - \alpha \mathbf{A} \mathbf{D}^{-1})^{-1} \mathbf{1}$, and thus, as before, β plays the role only of an unimportant overall multiplier

for the centrality. Conventionally, we set $\beta = 1$, giving

$$\mathbf{x} = (\mathbf{I} - \alpha \mathbf{A} \mathbf{D}^{-1})^{-1} \mathbf{1}. \quad (7.11)$$

This centrality measure is commonly known as *PageRank*, which is a name given it by the Google web search corporation. Google uses PageRank as a central part of their web ranking technology for web searches, which estimates the importance of web pages and hence allows the search engine to list the most important pages first [82]. PageRank works for the Web precisely because having links to your page from important other pages is a good indication that your page may be important too. But the added ingredient of dividing by the out-degrees of pages ensures that pages that simply point to an enormous number of others do not pass much centrality on to any of them, so that, for instance, network hubs like Amazon do not have a disproportionate influence on the rankings. PageRank also finds applications in other areas besides web search—see Gleich [205] for a review.

As with the Katz centrality, the formula for PageRank, Eq. (7.11), contains a free parameter α , whose value must be chosen somehow before the algorithm can be used. By analogy with Eq. (7.8) and the argument that follows it, we can see that the value of α should be less than the inverse of the largest eigenvalue of $\mathbf{A} \mathbf{D}^{-1}$. For an undirected network this largest eigenvalue turns out to be one,⁸ and thus α should be less than one. There is no equivalent result for a directed network, the leading eigenvalue differing from one network to another, although it is usually still roughly of order one.

The Google search engine uses a value of $\alpha = 0.85$ in its calculations, although it's not clear that there is any rigorous theory behind this choice. More likely it is just a shrewd guess based on experimentation to find out what works.

One could imagine a version of the PageRank equation (7.9) that did not have the additive constant term β in it at all:

$$x_i = \alpha \sum_j A_{ij} \frac{x_j}{k_j}, \quad (7.12)$$

which is similar to the original eigenvector centrality introduced back in Section 7.1.2, but now with the extra division by k_j . Particularly for undirected

⁸This is straightforward to show. The corresponding (right) eigenvector is (k_1, k_2, k_3, \dots) , where k_i is the degree of the i th node. It is easily confirmed that this is indeed an eigenvalue of $\mathbf{A} \mathbf{D}^{-1}$ with eigenvalue 1. Moreover, since this vector has all elements non-negative it must be the leading eigenvector and 1 must be the largest (most positive) eigenvalue by the Perron–Frobenius theorem discussed in Section 7.1.2—see footnote 2 on page 161.

	With constant term	Without constant term
Divide by out-degree	$\mathbf{x} = (\mathbf{I} - \alpha \mathbf{A} \mathbf{D}^{-1})^{-1} \mathbf{1}$ PageRank	$\mathbf{x} = \mathbf{A} \mathbf{D}^{-1} \mathbf{x}$ degree centrality
No division	$\mathbf{x} = (\mathbf{I} - \alpha \mathbf{A})^{-1} \mathbf{1}$ Katz centrality	$\mathbf{x} = \kappa^{-1} \mathbf{A} \mathbf{x}$ eigenvector centrality

Table 7.1: Four centrality measures. The four matrix-based centrality measures discussed in the text are distinguished by whether they include an additive constant term in their definition and whether they are normalized by dividing by node degrees. The matrix \mathbf{D} is the diagonal matrix with elements $D_{ii} = \max(k_i, 1)$ for undirected networks or $\max(k_i^{\text{out}}, 1)$ for directed ones—see Eq. (7.9) and the following discussion. Each of the measures can be applied to directed networks as well as undirected ones, although only PageRank and Katz centrality are commonly used in this way. The measure that appears in the top right corner of the table is equivalent to degree centrality in the undirected case. It takes more complicated values in the directed case but is not widely used.

networks, where the added β term is not really needed, this definition might appear to make sense. In fact, however, it turns out to be trivial for the undirected case: it is easy to see that Eq. (7.12) has solution $x_i = k_i$ (and $\alpha = 1$) on an undirected network and therefore is just the same as ordinary degree centrality. For a directed network it does not reduce to any equivalent simple value and it might potentially be of use, but then it suffers from the same problem as the original eigenvector centrality, in that it gives non-zero scores only to nodes that fall in strongly connected components of two or more nodes or their out-components. All other nodes get a zero score. Overall, therefore, this measure is not ideal and it does not find much practical use.

In Table 7.1 we give a summary of the different matrix centrality measures we have discussed, organized according to their definitions and properties. If you want to use one of these measures in your own calculations and find the alternatives bewildering, eigenvector centrality and PageRank are probably the two measures to focus on initially, being the most commonly used. The Katz centrality has found use in the past but has been favored less in recent work, while the PageRank measure without the constant term, Eq. (7.12), is the same as degree centrality for undirected networks and not in common use for directed ones.

7.1.5 HUBS AND AUTHORITIES

The centrality measures we have considered so far for directed networks all follow basically the same principle: they accord a node high centrality if it is pointed to by others with high centrality. However, in some networks it is appropriate also to accord a node high centrality if it *points to* others with high centrality. For instance, in citation networks there are review articles that cite a selection of notable papers on a certain subject. A review may itself contain relatively little information on the subject in question, but it tells us where to find information, and this on its own makes the review useful. Similarly, there are web pages that consist primarily of links to other pages on a given topic or topics and such a page of links could be very useful even if it does not itself contain explicit information on the topic in question.

Thus there are really two types of important nodes in these networks: *authorities* are nodes that contain useful information on a topic of interest and *hubs* are nodes that tell us where the best authorities are to be found.⁹ An authority may also be a hub and vice versa: review articles, for instance, do often contain useful discussions of the topic at hand as well as citations to other discussions. Clearly hubs and authorities only exist in directed networks, since in the undirected case there is no distinction between pointing to a node and being pointed to.

The concept of hubs and authorities in networks was first put forward by Kleinberg [265] and developed by him into a centrality algorithm called *hyperlink-induced topic search* or *HITS*. The HITS algorithm gives each node i in a directed network two different centrality scores, the *authority centrality* x_i and the *hub centrality* y_i , which quantify nodes' prominence in the two roles. The defining characteristic of a node with high authority centrality is that it is pointed to by many nodes with high hub centrality. Conversely, the defining characteristic of a node with high hub centrality is that it *points to* many nodes with high authority centrality.

Thus an important scientific paper (in the authority sense) would be one cited in many important reviews (in the hub sense). An important review is one that cites many important papers. But ordinary papers can also have high hub centrality if they cite many other important papers, and reviews may be cited by other hubs and hence have high authority centrality.

⁹In Chapter 1 we used the word "hub" in a different sense to mean a node with particularly high degree (see also Section 10.3). Both uses of the word are common in the networks literature, which can be confusing. When talking about hubs in this book we will be careful to make clear which sense we have in mind, and you should do the same.

In Kleinberg's approach the authority centrality of a node is defined to be proportional to the sum of the hub centralities of the nodes that point to it:

$$x_i = \alpha \sum_j A_{ij} y_j, \quad (7.13)$$

where α is a constant. Similarly, the hub centrality of a node is proportional to the sum of the authority centralities of the nodes it points to:

$$y_i = \beta \sum_j A_{ji} x_j, \quad (7.14)$$

with β another constant. Note that the indices on the matrix element A_{ji} are swapped around in this second equation: it is the nodes that i points to that define its hub centrality.

In matrix terms these equations can be written as

$$\mathbf{x} = \alpha \mathbf{A} \mathbf{y}, \quad (7.15)$$

$$\mathbf{y} = \beta \mathbf{A}^T \mathbf{x}, \quad (7.16)$$

or, combining the two,

$$\mathbf{A} \mathbf{A}^T \mathbf{x} = \lambda \mathbf{x}, \quad (7.17)$$

and

$$\mathbf{A}^T \mathbf{A} \mathbf{y} = \lambda \mathbf{y}, \quad (7.18)$$

where $\lambda = (\alpha\beta)^{-1}$. Thus the authority and hub centralities are respectively given by eigenvectors of $\mathbf{A} \mathbf{A}^T$ and $\mathbf{A}^T \mathbf{A}$ with the same eigenvalue. As with the standard eigenvector centrality of Section 7.1.1, assuming we want the centrality scores to be non-negative, the Perron–Frobenius theorem tells us that out of the n possible eigenvectors we must take the one corresponding to the largest (most positive) eigenvalue.

A crucial condition for this approach to work is that $\mathbf{A} \mathbf{A}^T$ and $\mathbf{A}^T \mathbf{A}$ have the same leading eigenvalue λ . Otherwise we cannot satisfy both Eq. (7.17) and Eq. (7.18). It is easily proved, however, that this is the case, and in fact that all eigenvalues are the same for the two matrices. If $\mathbf{A} \mathbf{A}^T$ has an eigenvalue λ , so that $\mathbf{A} \mathbf{A}^T \mathbf{x} = \lambda \mathbf{x}$ for some \mathbf{x} , then multiplying both sides by \mathbf{A}^T gives

$$\mathbf{A}^T \mathbf{A} (\mathbf{A}^T \mathbf{x}) = \lambda (\mathbf{A}^T \mathbf{x}), \quad (7.19)$$

which tells us that $\mathbf{A}^T \mathbf{A}$ also has an eigenvalue λ (with corresponding eigenvector $\mathbf{A}^T \mathbf{x}$).

Note that there is no need in practice to solve both the eigenvalue equations (7.17) and (7.18): if we solve for \mathbf{x} then we can calculate \mathbf{y} from Eq. (7.16).

(The factor β is unknown, but it multiplies all elements of \mathbf{y} equally and so does not affect their relative values, only their overall magnitude. Usually we do not care about overall magnitude, but if we do then it can be fixed for both \mathbf{x} and \mathbf{y} by normalizing in any convenient fashion.)

A nice feature of the hub and authority centralities is that they circumvent the problems that ordinary eigenvector centrality has with directed networks, that only nodes in strongly connected components of size two or more, or their out-components, have non-zero centrality. In the hub and authority approach nodes not cited by any others have authority centrality zero (which is reasonable), but they can still have non-zero hub centrality. And the nodes that *they* cite can then have non-zero authority centrality by virtue of being cited. This is perhaps a more elegant solution to the problems of eigenvector centrality in directed networks than the more ad hoc method of introducing an additive constant term as we did in Eq. (7.5). (We can employ such a constant term in the HITS algorithm if we wish, although there seems little point. We can also apply any of the other approaches considered in previous sections such as dividing node centralities by their out-degrees. Some variations along these lines are explored in Refs. [80,372], but we leave the pursuit of such details to the enthusiastic reader.)

The HITS algorithm is an elegant construct that should in theory provide more information about node centrality than the simpler measures of previous sections. In practice, however, it has not found much application. It was at one time used as the basis for the (now defunct) web search engines *Teoma* and *Ask Jeeves*, and will perhaps in future find other uses, for example in citation networks, where it has advantages over other eigenvector measures.

7.1.6 CLOSENESS CENTRALITY

In the preceding sections we have examined a number of centrality measures based on matrix concepts, particularly eigenvectors, but this is by no means the only formulation of centrality. In this section and the following one we look at two entirely different measures of centrality, both based on shortest paths in networks.

Closeness centrality is a centrality score that measures the mean distance from a node to other nodes. In Section 6.11.1 we encountered the concept of the shortest distance through a network between two nodes, i.e., the number of steps along the shortest path. Suppose d_{ij} is the shortest distance from node i to

Recall that shortest paths need not be unique—nodes can be joined by several shortest paths of the same length. The length d_{ij} , however, is always well defined, being the length of any one of those paths.

node j . Then the mean shortest distance from i to every node in the network is¹⁰

$$\ell_i = \frac{1}{n} \sum_j d_{ij}. \quad (7.20)$$

This quantity takes low values for nodes that are separated from others by only a short distance on average. It is plausible that such nodes might have more direct influence on others or better access. In a social network, for instance, a person with lower mean distance to others might find that their opinions spread through the community more quickly than other people's.

The mean distance ℓ_i is not a centrality measure in the same sense as the others in this chapter, since it gives low values for more central nodes and high values for less central ones, which is the opposite of our other measures. In the social networks literature, therefore, researchers commonly calculate the inverse of ℓ_i rather than ℓ_i itself. This inverse is called the *closeness centrality* C_i :

$$C_i = \frac{1}{\ell_i} = \frac{n}{\sum_j d_{ij}}. \quad (7.21)$$

Closeness centrality is a very natural measure of centrality, often used in social and other network studies. For example, it has become popular in recent years to rank film actors according to their closeness centrality in the network of who has appeared in films with whom [466]. Using data from the Internet Movie Database,¹¹ we find that in the largest component of the network, which includes more than 98% of all actors, the highest closeness centrality is 0.4143 for the actor Christopher Lee. This is an interesting result: Lee was neither as famous nor as successful as some of his contemporaries. He did, however, appear in an extraordinary number of films over the course of his long career—over 200 according to the database. This alone has a tendency to reduce his average distance to other nodes in the network by increasing the number of his

Lee is perhaps best known for his role as the wizard Saruman in the film versions of *The Lord of the Rings* and *The Hobbit*.

¹⁰In calculating the average distance some authors exclude from the sum in (7.20) the term for $j = i$, so that

$$\ell_i = \frac{1}{n-1} \sum_{j(\neq i)} d_{ij},$$

which is a reasonable strategy, since a node's influence on itself is usually not relevant to the working of the network. On the other hand, the distance d_{ii} from i to itself is zero by definition, so this term in fact contributes nothing to the sum. The only difference the change makes to ℓ_i is in the leading divisor, which becomes $1/(n-1)$ instead of $1/n$. Since this change is independent of i and since, as we have said, we usually care only about the relative centralities of different nodes and not their absolute values, we can in most cases ignore the difference between the two definitions. In this book we use (7.20) because it tends to give slightly more elegant analytic results.

¹¹<http://www.imdb.com>

collaborators and hence creating more paths through the network. Indeed we should, in general, expect nodes with higher degree to have shorter average distance to others, meaning that closeness centrality and degree centrality are positively correlated.

One potential problem with the definition of closeness centrality in Eq. (7.21) concerns networks that have more than one component. If, as discussed in Section 6.11.1, we define the shortest distance between two nodes to be infinite if the nodes fall in different components, then ℓ_i is infinite for all i in any network with more than one component (because the sum has at least one term that is infinite) and C_i is zero. There are two possible strategies for getting around this problem. The most common one is simply to average over only those nodes in the same component as i . Then n in Eq. (7.21) becomes the number of nodes in the component and the sum is over only that component. This gives us a finite measure, but one that has its own problems. In particular, distances tend to be smaller between nodes in small components, so that nodes in such components get lower values of ℓ_i and higher closeness centrality than their counterparts in larger components. This is usually undesirable: in most cases nodes in small components are considered *less* well connected than those in larger ones and should therefore be given lower centrality.

Perhaps a better solution is to redefine closeness in terms of the harmonic mean distance between nodes, i.e., the average of the inverse distances:

$$C'_i = \frac{1}{n-1} \sum_{j(\neq i)} \frac{1}{d_{ij}}. \quad (7.22)$$

(Note that we are obliged in this case to exclude from the sum the term for $j = i$, since $d_{ii} = 0$ which would make this term infinite. This means that the sum has only $n - 1$ terms in it, hence the leading factor of $1/(n - 1)$.)

This definition has a couple of nice properties. First, if $d_{ij} = \infty$ because i and j are in different components, then the corresponding term in the sum is simply zero and drops out. Second, the measure naturally gives more weight to nodes that are close to i than to those far away. Intuitively we might imagine that the distance to close nodes is what matters most in practical situations—once a node is far away in a network it matters less exactly how far away it is, and Eq. (7.22) reflects this, having contributions close to zero from all such nodes.

Despite its desirable qualities, however, Eq. (7.22) is rarely used in practice. The author has seen it employed only occasionally.

7.1.7 BETWEENNESS CENTRALITY

A different concept of node importance is captured by *betweenness centrality*, which measures the extent to which a node lies on paths between other nodes. The idea of betweenness is usually attributed to Freeman in 1977 [189], although as Freeman himself has pointed out [191] it was independently proposed some years earlier by Anthonisse in an unpublished technical report [23].

Suppose we have a network with something flowing around it from node to node along the edges. On the Internet, for instance, we have data packets flowing around. In a social network we might have messages, news, information, or rumors being passed from one person to another. Let us make the simple assumption that every pair of nodes in the network exchanges messages at the same average rate (more precisely every pair that is actually connected by a path) and that messages always take the shortest available path through the network, or one such path chosen at random if there are several. Then we ask the following question: if we wait a suitably long time until many messages have passed between every pair of nodes, how many messages will have passed through each node en route to their destination? The answer is that, since messages travel along shortest paths, the number passing through each node is proportional to the number of shortest paths the node lies on. This number of shortest paths is what we call the betweenness centrality, or just betweenness for short.

It seems plausible that nodes with high betweenness centrality could have influence within a network, by virtue of their control over information passing between others. The nodes with highest betweenness in our message-passing scenario are the ones through which the largest number of messages pass, and if those nodes get to see the messages in question as they go by, or if nodes get paid for passing messages along, then they could derive a lot of power or wealth from their position within the network. The nodes with highest betweenness are also the ones whose removal from the network will most disrupt communications between other nodes in the sense that they lie on the largest number of paths taken by messages. If a node with high betweenness is removed then all the messages that would have passed through that node must now be rerouted another way.

In real-world situations it is usually not the case that all nodes exchange messages with the same frequency or that messages always take the shortest path. Nonetheless, betweenness centrality may still be a reasonable guide to the influence nodes have over the flow of information between others.

Mathematically, betweenness centrality can be expressed as follows. Suppose for the moment that we have an undirected network in which there is at

most one shortest path between any pair of nodes. (There may be zero paths if the nodes in question are in different components.) Let n_{st}^i be 1 if node i lies on the shortest path from s to t and 0 if it does not or if there is no such path. Then the betweenness centrality x_i is given by

$$x_i = \sum_{st} n_{st}^i. \quad (7.23)$$

Note that this definition counts separately the shortest paths in either direction between each node pair. On an undirected network these paths are the same, so this effectively counts each path twice. One could compensate for this by dividing x_i by 2, and sometimes this is done, but we prefer the definition given here for a couple of reasons. First, it makes little difference in practice whether one divides the centrality by 2, since one is usually concerned only with the relative magnitudes of the centralities and not with their absolute values. Second, as discussed below, Eq. (7.23) has the advantage that it can be applied unmodified to directed networks, in which the paths in either direction between a node pair can differ.

Note also that Eq. (7.23) includes paths from each node to itself. Some people prefer to exclude such paths from the definition, so that $x_i = \sum_{s \neq t} n_{st}^i$, but again the difference is typically not important. Every node lies on one path from itself to itself, so the inclusion of these terms simply increases the betweenness by 1, but does not change the rankings of the nodes—which ones have higher or lower betweenness—relative to one another.¹²

Equation (7.23) applies in the case where there is at most one shortest path between each pair of nodes. More generally there may be more than one (see Section 6.11.1). The standard extension of betweenness to this case gives each path between two nodes a weight equal to the inverse of the number of paths, so that, for instance, if there are two shortest paths between a given pair of nodes, each of them gets weight $\frac{1}{2}$. Then the betweenness of a node is defined to be the sum of the weights of all shortest paths passing through that node.

¹²One could also ask whether the path from s to t should be considered to pass through the nodes s and t themselves. In the social networks literature it is usually assumed that it does not, but we prefer the definition given here where it does: it seems reasonable to define a node to be on a path between itself and someone else, since normally a node has control over information flowing from itself to other nodes or vice versa. If, however, we exclude the endpoints of the path, the only effect is to reduce the number of paths through each node by twice the size of the component to which the node belongs. Thus the betweennesses of all nodes within a single component are reduced by the same additive constant and the ranking of nodes within the component is again unchanged. The rankings of nodes in different components can change relative to one another, but this is rarely an issue because betweenness centrality is not typically used to compare nodes in different components, since such nodes are not competing for influence in the same arena.

Note that it is possible for two shortest paths between the same pair of nodes to overlap, sharing some nodes in addition to the starting and ending nodes—see Fig. 7.2. If two or more paths pass through the same node then the betweenness sum includes contributions from each of them.

Formally, we redefine n_{st}^i to be the number of shortest paths from s to t that pass through i and we define g_{st} to be the total number of shortest paths from s to t . Then the betweenness centrality of node i on a general network is

$$x_i = \sum_{st} \frac{n_{st}^i}{g_{st}}, \quad (7.24)$$

where we adopt the convention that $n_{st}^i/g_{st} = 0$ if both n_{st}^i and g_{st} are zero. This definition is equivalent to our message-passing thought experiment above, in which messages travel along shortest paths between nodes and when there is more than one shortest path they choose a path at random. Then the betweenness of Eq. (7.24) is proportional to the average rate at which traffic passes through node i .

So far we have considered only undirected networks, but betweenness centrality can be applied to directed networks as well. In a directed network the shortest path between two nodes depends, in general, on the direction you travel in. The shortest path from A to B is different from the shortest path from B to A . Indeed there may be a path in one direction and no path at all in the other. Thus it is important in a directed network explicitly to include the path counts in either direction between each node pair. The definition in Eq. (7.24) already does this, so we can use the same definition without modification for the directed case. Although the generalization of betweenness to directed networks is straightforward, however, it is rarely used and we will not consider it further in this book.

Betweenness centrality differs from the other centrality measures we have considered in being not principally a measure of how well-connected a node is. Instead it measures how much a node falls “between” others. Indeed a node can have quite low degree, be connected to others that have low degree, even be a long way from others on average, and still have high betweenness. Consider the situation depicted in Fig. 7.3. Node A lies on a bridge between two groups within the network. Since any shortest path (or indeed any path whatsoever) between a node in group 1 and a node in group 2 must pass along this bridge, A has high betweenness centrality, even though it is itself on the

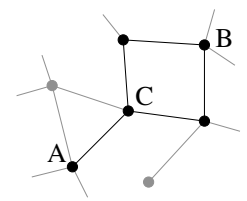


Figure 7.2: Overlapping shortest paths. Nodes A and B in this network are connected by two shortest paths, and node C lies on both paths.

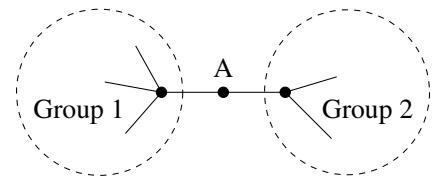


Figure 7.3: A low-degree node with high betweenness. In this sketch of a network, node A lies on a bridge joining two groups of other nodes. All paths between the groups must pass through A , so it has a high betweenness even though its degree is low.

periphery of both groups and has low degree. Probably it would not have particularly impressive values for eigenvector or closeness centrality, and yet it might have a lot of influence in the network as a result of its control over the flow of information between others. Nodes in roles like this are sometimes referred to in the sociological literature as *brokers*.¹³

The second highest betweenness score goes to Christopher Lee again.

As an example of betweenness centrality, consider again the network of film actors from the previous section. It turns out that the individual with highest betweenness in this network is the great Spanish actor Fernando Rey, most famous in the English-speaking world for his 1971 starring role next to Gene Hackman in *The French Connection*. It is perhaps no coincidence that the highest betweenness belongs to an actor who appeared in both European and American films, played roles in several different languages, and worked extensively in both film and television, as well as on stage. Rey was the archetypal broker, a crucial link between different branches of the entertainment industry.

The values of betweenness as defined above are raw numbers of paths, but it is sometimes convenient to normalize betweenness in some way. (Several of the standard computer programs for network analysis, such as Pajek and UCINET, perform such normalizations.) One natural choice is to normalize the path count by dividing by the total number of (ordered) node pairs, which is n^2 , so that betweenness becomes the fraction (rather than the number) of paths that run through a given node:¹⁴

$$x_i = \frac{1}{n^2} \sum_{st} \frac{n_{st}^i}{g_{st}}. \quad (7.25)$$

With this definition, the values of the betweenness lie strictly between zero and one.

A number of variations on the betweenness centrality have been proposed, mainly aimed at expanding the set of paths considered beyond just shortest paths, since real network traffic often doesn't take the shortest path to its destination. Many of us, for instance, have had the experience of hearing news

¹³Much sociological literature addresses concepts of power or "social capital." It may seem ruthless to think of individuals exploiting their control over other people's information to gain the upper hand, but it may also be realistic. At least in situations where there is a significant pay-off to having such an upper hand (like business relationships, for example), it is reasonable to suppose that notions of power derived from network structure really do figure in people's interactions with the world around them.

¹⁴Another possibility, proposed by Freeman in his original paper on betweenness [189], is to divide by the maximum value that betweenness can take on any network of size n , which for our definition of betweenness is $n^2 - n + 1$. We prefer Eq. (7.25) for its ease of interpretation, although the difference between the two becomes small anyway in the limit of large n .

about one of our friends not from that friend directly but from another mutual acquaintance—the message has passed along a path of length two via the mutual acquaintance, rather than along the direct (shortest) path of length one.

Flow betweenness is a variant of betweenness centrality that uses edge-independent paths between node pairs rather than shortest paths [192]. If there is more than one possible choice of independent paths between a pair of nodes, the contribution to the betweenness of any node for that pair is defined to be the maximum over all choices.

Another variant is *random-walk betweenness* [356], which imagines messages performing random walks across the network between every possible starting point and destination, and the betweenness is defined as the average number of such messages that pass through each node. Random-walk betweenness would be an appropriate betweenness measure for traffic that traverses a network with no idea of where it is going—it simply wanders around at random until it reaches its destination. Conventional shortest-path betweenness is the exact opposite: it is the appropriate measure for information that knows exactly where it is going and takes the most direct route to get there. It seems likely that most real-world situations fall somewhere in between these two extremes. It is found in practice, however, that the two measures often give quite similar results [356], in which case one can with reasonable justification assume that the “correct” answer, which presumably lies between the limits set by the shortest-path and random-walk measures, is similar to both. In cases where the two differ by a larger margin, however, we should be wary of attributing too much authority to either measure—there is no guarantee that either is telling us a great deal about true information flow in the network.

Other generalizations of betweenness are also possible, based on other models of diffusion, transmission, or flow along network edges. We refer the interested reader to the article by Borgatti [76], which draws together many of the possibilities into a broad general framework for betweenness measures.

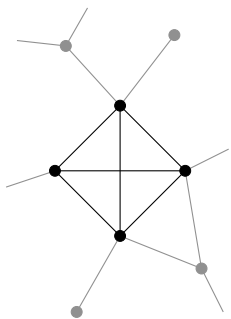
7.2 GROUPS OF NODES

Many networks, including social and other networks, divide naturally into groups or communities. Networks of people divide into groups of friends, co-workers, or business partners; the World Wide Web divides into groups of related web pages; biochemical networks divide into functional modules, and so forth. The definition and analysis of groups within networks is a large and fruitful area of network theory. In Chapter 14 we discuss some of the sophisticated computer methods that have been developed for dividing networks into their constituent groups, such as modularity-based methods

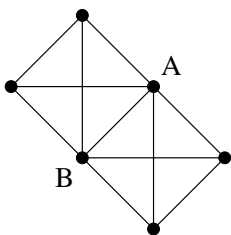
See Section 6.13 for a discussion of independent paths.

See Section 6.14.3 for a discussion of random walks.

and maximum likelihood methods. In this section we discuss some simpler concepts of network groups that can be useful for probing and describing the local structure of networks. The primary constructs we look at are cliques, k -cores, and k -components.



A clique of four nodes within a network.



Two overlapping cliques. Nodes A and B in this network both belong to two cliques of four nodes.

7.2.1 CLIQUES

A *clique* is a set of nodes within an undirected network such that every member of the set is connected by an edge to every other. Thus a set of four nodes in a network would be a clique if (and only if) each of the four is directly connected by edges to the other three. Note that cliques can overlap, meaning that they can share one or more of the same nodes.

The occurrence of a clique in an otherwise sparsely connected network is normally an indication of a highly cohesive subgroup. In a social network, for instance, one might encounter a set of individuals each of whom was acquainted with each of the others, and such a clique would probably indicate that the individuals in question are closely connected—the members of a family, for example, or a set of co-workers in an office.

However, it's also the case that many circles of acquaintances form only near-cliques, rather than perfect cliques. There may be some members of a group who are unacquainted, even if most members know one another. The requirement that every possible edge be present within a clique is a very stringent one and limits the usefulness of the clique concept. There are, however, some circumstances in which cliques do crop up and play an important role. An example is the one-mode projection of a bipartite network introduced in Section 6.6.1. Recall that bipartite networks (also called affiliation networks in sociology) are commonly used to represent the membership of people or objects in groups of some kind. The one-mode projection creates a network that is naturally composed of cliques, one for each group—see Fig. 6.6 on page 117.

7.2.2 CORES

For many purposes a clique is too stringent a notion of grouping to be useful and it is natural to ask how one might define something more flexible. One possibility is the k -core. By contrast with a clique, where each node is joined to all the others, a k -core is a connected set of nodes where each is joined to at least k of the others. Thus, in a 2-core, for instance, every node is joined to at least two others in the set. Figure 7.4 shows the k -cores in a small network.

The k -core is not the only possible relaxation of a clique, but it is a particularly useful one for the very practical reason that k -cores are easy to find.

Note that a 1-core is the same thing as an ordinary component.

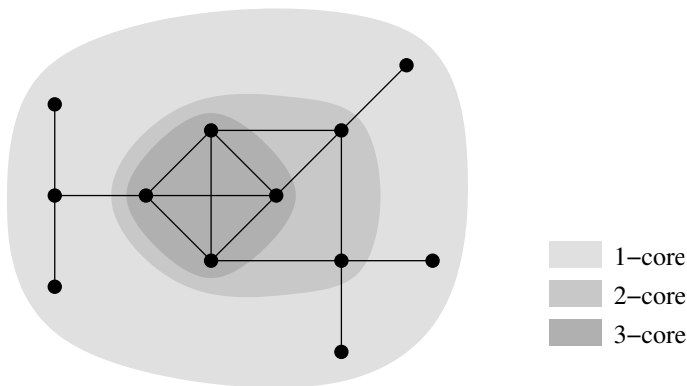


Figure 7.4: The k -cores in a small network. The shaded regions denote the k -cores for $k = 1, 2$, and 3 in this small network. There are no k -cores for $k > 3$ in this case. Note how the k -cores are nested within one another, the 3-core inside the 2-core, which is in turn inside the 1-core.

A simple way to find them is to start with a given network and remove from it any nodes that have degree less than k , along with their attached edges, since clearly such nodes cannot under any circumstances be members of a k -core. In so doing, one will normally reduce the degrees of some other nodes in the network—those that were connected to the nodes just removed. So we then go through the network again to see if there are any additional nodes that now have degree less than k and remove those too. And so we proceed, repeatedly pruning the network to remove nodes with degree less than k until no such nodes remain. What is left over will, by definition, be a k -core or a set of k -cores, since each node is connected to at least k others. Note that we are not necessarily left with a *single* k -core—there’s no guarantee that the network will be connected once we are done pruning it, even if it was connected to start with.

For any given network, there is a maximum value of k for the k -cores. It is clear, for instance, that no k -cores can exist when k exceeds the highest degree in the network, since in that case no node could have k connections to others. The k -cores of a network also have the property of being *nested* within one another: the 2-cores are subsets of the 1-cores, the 3-cores subsets of the 2-cores, and so forth—see Fig. 7.4. This must be the case since one could, if one wished, compute the 3-cores by first removing all nodes with degree less than 2, thereby creating the 2-cores, then removing all nodes with degree less than 3 from those, creating the 3-cores. Thus, the breakdown of a network

There is a close connection between k -cores and the concept of “complex contagion,” which is used to model the spread of ideas or information in social networks. See the discussion in Sections 16.1.9 and 16.3.5 and footnote 12 on page 640. Another closely related process, bootstrap percolation, has been studied extensively in statistical physics—see Refs. [7, 99, 210].

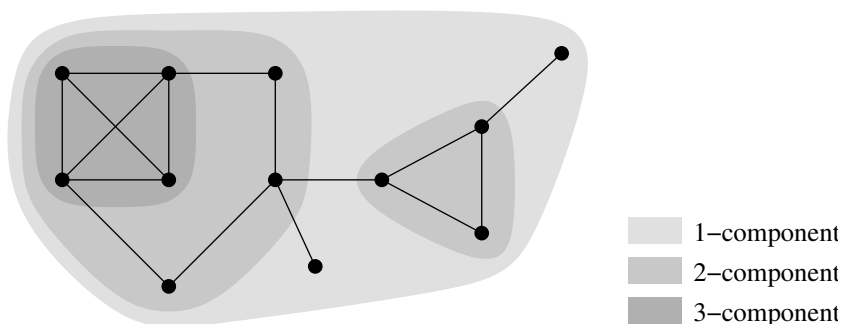


Figure 7.5: The k -components in a small network. The shaded regions denote the k -components in this small network, which has a single 1-component, two 2-components, one 3-component, and no k -components for any higher value of k . Note that the k -components are nested within one another, the 2-components falling inside the 1-component and the 3-component falling inside one of the 2-components.

into cores for all values of k provides a onion-like decomposition into layers within layers—1-cores, then 2-cores, then 3-cores, and so forth, culminating at the highest value of k for which cores exist. This decomposition is sometimes used as a measure of *core-periphery structure* in networks: nodes that lie within the highest- k cores are “core” nodes within the network, while nodes outside those cores are “peripheral” nodes. In this sense, the cores define a kind of centrality measure, and they are sometimes used that way. In the social networks literature, for instance, it is sometimes hypothesized that core actors in a network, defined in this sense, may be more powerful or influential, or have better access to information or resources, although this is only a hypothesis—there is in most cases no formal reason to suppose that k -cores are closely linked with node roles or behaviors [462].

7.2.3 COMPONENTS AND k -COMPONENTS

In Section 6.12 we introduced the concept of a component. A component in an undirected network is a (maximal) set of nodes such that each is reachable by some path from each of the others. A useful generalization of this concept is the k -component. A k -component (sometimes also called a k -connected component) is a set of nodes such that each is reachable from each of the others by at least k node-independent paths—see Fig. 7.5. (Recall that two paths are said to be node-independent if they share none of the same nodes except the starting and ending nodes—see Section 6.13.) For the common special cases $k = 2$ and $k = 3$,

See Section 14.7.3 for further discussion of core-periphery structure.

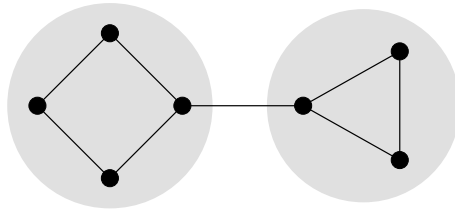


Figure 7.6: A small network with one 2-core but two 2-components. The whole of this network constitutes a single 2-core, since each of its nodes is connected to at least two of the others. But the network contains two separate 2-components, as indicated by the two shaded circles, proving that 2-cores and 2-components are not the same thing.

k -components are also called *bicomponents* and *tricomponents* respectively.

A 1-component by this definition is just an ordinary component—there is at least one path between every pair of nodes—and, like the k -cores of the previous section, k -components are nested within each other. A 2-component or bicomponent, for example, is necessarily a subset of a 1-component, since any pair of nodes that are connected by at least two paths are also connected by at least one path. Similarly a tricomponent is necessarily a subset of a bicomponent, and so forth. (See Fig. 7.5 again.)

At first sight, k -components seem rather similar to k -cores, but there are important differences. Consider Fig. 7.6, which shows a small network which is composed of a single 2-core—every node in the network is connected to at least two of the others—yet there are two separate 2-components in the network. The left and right halves of the network are connected by only one independent path in the middle, so they are separate 2-components.

As discussed in Section 6.13, the number of node-independent paths between two nodes is equal to the size of the minimum node cut set between the same two nodes, i.e., the number of nodes that would have to be removed in order to disconnect the two. So another way of defining a k -component would be to say that it is a subset of a network in which no pair of nodes can be disconnected from each other by removing less than k other nodes.

A variant of the k -component can also be defined using edge-independent paths, so that nodes are in the same k -component if they are connected by k or more edge-independent paths, or equivalently if they cannot be disconnected by the removal of less than k edges. In principle this variant could be useful in certain circumstances but in practice it is rarely used.

The idea of a k -component is a natural one in network analysis, being connected with the idea of network robustness. For instance, in a data network

such as the Internet, the number of node-independent paths between two nodes is also the number of independent routes that data might take between the same two nodes, and the size of the cut set between them is the number of nodes in the network—typically routers—that would have to fail or otherwise be knocked out to sever the data connection between the two endpoints. Thus a pair of nodes connected by two independent paths cannot be disconnected from one another by the failure of any single router. A pair of nodes connected by three paths cannot be disconnected by the failure of any two routers. And so forth. A k -component with $k \geq 2$ in a network like the Internet is a subset of the network that has robust connectivity in this sense. One would hope, for instance, that most of the network backbone—the system of high volume world-spanning links that carry long-distance data (see Section 2.1)—is a k -component with high k , so that it would be difficult for points on the backbone to lose connection with one another.

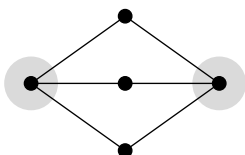


Figure 7.7: A non-contiguous tricomponent. The two highlighted nodes in this network form a tricomponent, even though they are not directly connected to each other. The other three nodes are not in the tricomponent.

One disadvantage of k -components as a definition of node groups, is that for $k \geq 3$ they can be non-contiguous (see Fig. 7.7). Ordinary components (1-components) and 2-components are always contiguous, but 3-components and above may not be. Within the social networks literature, where non-contiguous components are often considered undesirable, k -components are sometimes defined slightly differently, to be a set of nodes such that every pair in the set is connected by at least k node-independent paths *that themselves are contained entirely within the subset*. This definition rules out non-contiguous k -components, but it is also mathematically and computationally more difficult to work with than the standard definition. For this reason, and because there are also plenty of cases in which it is appropriate to count non-contiguous k -components, the standard definition remains the one most widely used.

There are a number of other definitions of node groups that find occasional use, particularly in the social networks literature, such as k -plexes and k -clubs. See the book by Wasserman and Faust [462] for a detailed discussion. There are also various definitions that avoid the use of a parameter k . For instance, Flake *et al.* [181] proposed a definition of a group as a set of nodes that each has at least as many connections inside the set as outside. Radicchi *et al.* [395] proposed a weaker definition where a group is a set of nodes such that the total number of connections between nodes inside the set is greater than the total number to nodes outside it. The use of these measures is, however, relatively rare and we will not consider them further here.

7.3 TRANSITIVITY AND THE CLUSTERING COEFFICIENT

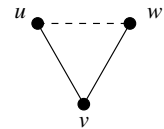
A notion particularly important in social networks, and useful to some degree in other networks too, is *transitivity*. In mathematics a relation “ \circ ” is said to be transitive if $a \circ b$ and $b \circ c$ together imply $a \circ c$. An example would be equality. If $a = b$ and $b = c$, then it follows also that $a = c$, so “ $=$ ” is a transitive relation. Other examples are “greater than,” “less than,” and “implies.”

The fundamental type of relation between nodes in a network is “connected by an edge.” If the “connected by an edge” relation were transitive it would mean that if node u is connected to node v , and v is connected to w , then u is also connected to w . In common parlance, “the friend of my friend is also my friend.” This is what we mean by network transitivity. It can apply to either directed or undirected networks, but let us take the undirected case first, since it’s simpler.

Perfect transitivity only occurs in networks where every component is a clique, i.e., all nodes in a component are connected to all others.¹⁵ Perfect transitivity is therefore not a very useful concept for most networks. However, *partial* transitivity can be useful. In many networks, particularly social networks, the fact that u knows v and v knows w doesn’t *guarantee* that u knows w , but makes it much more likely. The friend of my friend is not necessarily my friend, but is far more likely to be my friend than some randomly chosen member of the population.

We can quantify the level of transitivity in a network as follows. If u knows v and v knows w , then we have a path uvw , two edges long, in the network. If u also knows w , we say that the path is *closed*—it forms a loop of length three, or a triangle, in the network. In social network jargon, u , v , and w are said to form a *closed triad*. We define the *clustering coefficient*¹⁶ to be the fraction of paths of length two in the network that are closed. That is, we count all paths of length two, and we count how many of them are closed, and we divide the

See Section 7.2 for a discussion of cliques.



The path uvw (solid edges) is said to be closed if the third edge directly from u to w is present (dashed line).

¹⁵To see this suppose we have a component that is perfectly transitive but not a clique, meaning that there is at least one pair of nodes u, w in the component that are not directly connected by an edge. Since u and w are in the same component they must therefore be connected by some path of length greater than one, $u, v_1, v_2, v_3, \dots, w$. Consider the first two links in this path. Since u is connected by an edge to v_1 and v_1 to v_2 it follows that u must be connected to v_2 if the network is perfectly transitive. Then consider the next two links. Since u is connected to v_2 and v_2 to v_3 it follows that u must be connected to v_3 . Repeating the argument all the way along the path, we can then see that u must be connected by an edge to w . But this violates the hypothesis that u and w are not directly connected. Hence no perfectly transitive components exist that are not cliques.

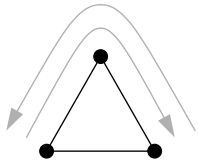
¹⁶The use of the word “clustering” in the name of the clustering coefficient is unconnected with the use of the same word in social network analysis to describe groups or clusters of nodes (see for instance Section 14.5.2). The reader should be careful to avoid confusing these two uses.

second number by the first to get a clustering coefficient C that lies in the range from zero to one:

$$C = \frac{(\text{number of closed paths of length two})}{(\text{number of paths of length two})}. \quad (7.26)$$

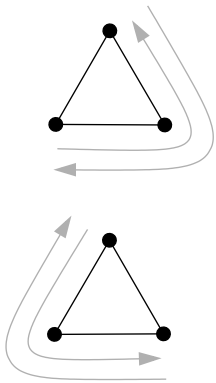
$C = 1$ implies perfect transitivity, i.e., a network whose components are all cliques. $C = 0$ implies no closed triads, which happens for various topologies, such as a tree (which has no closed loops of any kind—see Section 6.8) or a square lattice (which has closed loops with even numbers of nodes only but no closed triads).

Note that paths in networks, as defined in Section 6.11, have a direction (even in an undirected network). Thus uvw and wvu are considered distinct paths. The formula in Eq. (7.26) counts these paths separately, although in practice it would also be fine to count each path in only one direction—it would reduce both the numerator and the denominator by a factor of two, and the factors would cancel, leaving the value of C unchanged. Usually, however, and particular when writing computer programs, it is easier to count paths in both directions—it avoids having to remember which paths you have already counted.



An alternative way to write the clustering coefficient is

$$C = \frac{(\text{number of triangles}) \times 6}{(\text{number of paths of length two})}. \quad (7.27)$$



Why the factor of six? It arises because each triangle in a network contains six paths of length two. Suppose we have a triangle uvw . Then there are six paths of length two in it: uvw , vwu , wuv , wvu , vuw , and uwx . Each of these six is closed, so the number of closed paths is six times the number of triangles, and using this result in Eq. (7.26) then gives Eq. (7.27).

Another way to write the clustering coefficient would be to note that if we have a path of length two, uvw , then u and w have a common neighbor in v —they share a mutual acquaintance in social network terms. If the path uvw is closed then u and w are also themselves acquainted, so the clustering coefficient can be thought of also as the fraction of pairs of people with a common friend who are themselves friends, or equivalently as the mean probability that two people with a common friend are themselves friends.

This is perhaps the most common way of defining the clustering coefficient. In mathematical notation:

$$C = \frac{(\text{number of triangles}) \times 3}{(\text{number of connected triples})}. \quad (7.28)$$

A triangle contains six distinct paths of length two, all of them closed.

Here a “connected triple” means three nodes uvw with edges (u, v) and (v, w) . (The edge (u, w) can be present or not.) The factor of three in the numerator arises because each triangle gets counted three times when we count the connected triples in the network. The triangle uvw , for instance, contains the triples uvw , vuw , and wuv . In the older social networks literature the clustering coefficient is sometimes called the “fraction of transitive triples,” which is a reference to this definition of the coefficient.

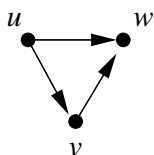
Social networks tend to have quite high values of the clustering coefficient. For example, the network of film actor collaborations discussed earlier in this chapter has $C = 0.20$ [354]; a network of collaborations between biologists was found to have $C = 0.09$ [349]; a network of who sends email to whom in a large university had $C = 0.16$ [156]. These are typical values for social networks. Some denser networks have even higher values, as high as 0.5 or 0.6. (Technological and biological networks by contrast tend to have somewhat lower values. The Internet at the autonomous system level, for instance, has a clustering coefficient of only about 0.01. This point is discussed in more detail in Section 10.6.)

In what sense are the clustering coefficients for social networks high? Let us assume, to make things simple, that everyone in a network has about the same number c of friends and let us suppose that everyone picks their friends completely at random from the whole population, meaning that they have the same probability of being friends with every person in the network. That probability is simply equal to $c/(n - 1)$, where n is the total number of people in the network. But in that case the probability of two of my friends being acquainted, which is by definition the clustering coefficient, is also $c/(n - 1)$ —my friends have the same probability of being acquainted as everyone else.

For the networks cited above, the value of $c/(n - 1)$ is 0.0003 (film actors), 0.00001 (biology collaborations), and 0.00002 (email messages). Thus the real clustering coefficients are *much* larger than our simple calculation would suggest. The calculation does ignore any variation in the number of friends people have, but the disparity between calculated and observed clustering coefficients is so large that it seems unlikely it could be eliminated just by allowing the number of friends to vary. A more likely explanation is that we were wrong to assume that everyone has the same probability of knowing everyone else. The numbers suggest that there is a much greater chance that two people will be acquainted if they have another common acquaintance than if they don’t. We discuss this point at greater length in Section 10.6.

Some social networks, such as the email network mentioned earlier, are directed networks. In calculating clustering coefficients for directed networks, scientists have typically just ignored their directed nature and applied Eq. (7.28)

Of course it is not normally the case that everyone in a network has the same number of friends. We will see later how to perform better calculations of the clustering coefficient (Section 12.3), but this simple calculation will serve our purposes for now.



A transitive triple of nodes in a directed network.

as if the edges were undirected. It is however possible to generalize transitivity to take account of directed links. If we have a directed relation between nodes such as “ u likes v ” then we can say that a triple of nodes is closed or transitive if u likes v , v likes w , and also u likes w . One can calculate a clustering coefficient in the obvious fashion for the directed case, counting all directed paths of length two that are closed and dividing by the total number of directed paths of length two. To date, however, such measurements have not often appeared in the literature.

7.3.1 LOCAL CLUSTERING AND REDUNDANCY

The clustering coefficient of the previous section is a property of an entire network. It quantifies the extent to which pairs of nodes with a common neighbor are also themselves neighbors, averaged over the whole network. It is, however, also sometimes useful to define a clustering coefficient for a single node. For a node i , we can define

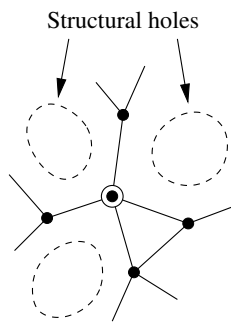
$$C_i = \frac{\text{(number of pairs of neighbors of } i \text{ that are connected)}}{\text{(number of pairs of neighbors of } i)}}. \quad (7.29)$$

In this book we use the notation C_i for both the local clustering coefficient and the closeness centrality. Care must be taken not to confuse the two.

That is, to calculate C_i we go through all distinct pairs of nodes that are neighbors of i , count the number of such pairs that are connected to each other, and divide by the total number of pairs, which is $\frac{1}{2}k_i(k_i - 1)$, where k_i is the degree of i . C_i is sometimes called the *local clustering coefficient* and it represents the average probability that a pair of i 's friends are friends of one another. (For nodes with degree zero or one the number of pairs of neighbors is zero and Eq. (7.29) is not well defined. Conventionally in this case we say that $C_i = 0$.)

Local clustering is interesting for several reasons. First, in many networks it is found empirically to have a rough dependence on degree, nodes with higher degree having a lower local clustering coefficient on average. This point is discussed in detail in Section 10.6.1.

Second, local clustering can be used as an indicator of so-called “structural holes” in a network. While it is common in many networks, especially social networks, for the neighbors of a node to be connected among themselves, it does happen sometimes that these expected connections between neighbors are missing. The missing links are called *structural holes* and were first studied in this context by Burt [89]. If we are interested in the efficient spread of information or other traffic around a network then structural holes are a bad thing—they reduce the number of alternative routes information can take through the network. On the other hand, structural holes could be a good thing for the node whose neighbors lack connections, because they give that



When the neighbors of a node are not connected to one another we say the network contains “structural holes.”

node power over information flow between those neighbors. If two of your neighbors are not connected directly and their information about one another comes via their mutual connection with you, then you can control the flow of that information. The local clustering coefficient measures how influential a node is in this sense, taking lower values the more structural holes there are in the surrounding network. Thus, local clustering can be regarded as a type of centrality measure, albeit one that takes small values for powerful individuals rather than large ones.

In this sense, local clustering can also be thought of as akin to the betweenness centrality of Section 7.1.7. Betweenness measures a node's control over information flowing between all pairs of nodes in its component. Local clustering is like a local version of betweenness that measures control over flows between just a node's immediate neighbors. One measure is not necessarily better than the other. There may be cases in which we want to take all nodes into account and others where we want to consider neighbors only. It is worth pointing out, however, that betweenness is much more computationally demanding to calculate than local clustering (see Section 8.5.6), and that in practice betweenness and local clustering are strongly correlated [89]. As a result, there may be in many cases little to be gained by performing the more costly full calculation of betweenness rather than using local clustering, given that the two contain much the same information.¹⁷

In his original studies of structural holes, Burt [89] did not make use of the local clustering coefficient,¹⁸ instead using another measure, which he called *redundancy*. The original definition of redundancy was rather complicated, but Borgatti [75] has shown that it can be simplified to the following: the redundancy R_i of a node i is the mean number of connections from a neighbor of i to other neighbors of i . Consider the example shown in Fig. 7.8 in which the central node has four neighbors. Each of those four *could* be acquainted with any of the three others, but in this case none of them is connected to all three. One is connected to none of the others, two are connected to one other, and the last is connected to two others. The redundancy is the average of these numbers $R_i = \frac{1}{4}(0 + 1 + 1 + 2) = 1$.

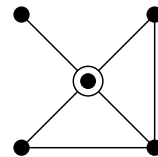


Figure 7.8: Redundancy. The neighbors of the central node in this figure have 0, 1, 1, and 2 connections to other neighbors, respectively. The redundancy is the mean of these values: $R_i = \frac{1}{4}(0 + 1 + 1 + 2) = 1$.

¹⁷As an example, in Section 14.5.1 we study methods for partitioning networks into clusters or communities and we will see that effective computer algorithms for this task can be created based on betweenness measures, but that almost equally effective and much faster algorithms can be created based on local clustering.

¹⁸Actually, the local clustering coefficient hadn't yet been invented. It was first proposed to this author's knowledge by Watts [463] a few years later.

The minimum possible value of the redundancy of a node i is zero and the maximum is $k_i - 1$, where k_i is the degree of the node.

Although redundancy and local clustering are different measures, they are related. To see what the relation is, we note that if the average number of connections from a friend of i to other friends is R_i , then the total number of connections between friends is $\frac{1}{2}k_i R_i$. And the total number of pairs of friends of i is $\frac{1}{2}k_i(k_i - 1)$. The local clustering coefficient, Eq. (7.29), is the ratio of these two quantities, so

$$C_i = \frac{\frac{1}{2}k_i R_i}{\frac{1}{2}k_i(k_i - 1)} = \frac{R_i}{k_i - 1}. \quad (7.30)$$

Given that $k_i - 1$ is the maximum value of R_i , the local clustering coefficient can thus be thought of as a version of the redundancy normalized to have a maximum value of 1. Applying Eq. (7.30) to the example of Fig. 7.8 implies that the local clustering coefficient for the central node should be $C_i = \frac{1}{3}$ and you can easily verify that this is indeed the case by calculating C_i directly from Eq. (7.29).

Another use of the local clustering coefficient is in the measurement of global clustering. Watts and Strogatz [466] proposed calculating a clustering coefficient for an entire network as the mean of the local clustering coefficients for each node,

$$C_{WS} = \frac{1}{n} \sum_{i=1}^n C_i, \quad (7.31)$$

where n is the number of nodes in the network. This is a different clustering coefficient from the one given earlier in Eq. (7.28)—the two are not equivalent—but both are in common use in the networks literature, which can lead to confusion. Furthermore, the two can in some cases give substantially different numbers when applied to the same network. We favor the first definition, Eq. (7.28), because it has a simple interpretation and because it is normally easier to calculate. Also the second definition, Eq. (7.31), can become dominated by nodes with low degree, since they have small denominators in Eq. (7.29), and the measure can thus give a poor picture of the overall properties of any network with a significant number of such nodes.¹⁹ Nonetheless, since both definitions are common you need to be aware of both and clear which is being used in any particular situation.

¹⁹As discussed in Section 10.6.1, nodes with low degree tend to have high values of C_i in most networks and this means that C_{WS} is usually larger than the value given by Eq. (7.28), sometimes much larger.

7.4 RECIPROCITY

The clustering coefficient of Section 7.3 measures the frequency with which triangles—loops of length three—appear in a network, but there is no reason why one should concentrate only on loops of length three and people have occasionally looked at the frequency of loops of length four or more [61, 92, 195, 207, 351]. Triangles occupy a special place, however, because in an undirected simple graph the triangle is the shortest loop we can have (and usually the most commonly occurring). However, in a *directed* network this is not the case. In a directed network, we can have loops of length two—a pair of nodes between which there are directed edges running in both directions—and it is interesting to ask about the frequency of occurrence of these loops also.

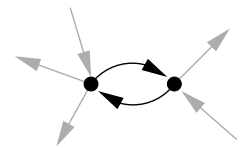
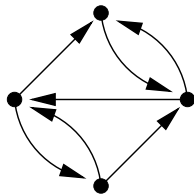
The frequency of loops of length two is measured by the *reciprocity*, which tells you how likely it is that a node you point to also points back at you. For instance, on the World Wide Web if my web page links to your web page, how likely is it, on average, that yours links back again to mine? In general, it's found that in fact you are much more likely to link to me if I link to you. Similarly, in friendship networks, such as those of Section 4.2, where respondents are asked to name their friends, it is much more likely that you will name me if I name you.

If there is a directed edge from node i to node j in a directed network and there is also an edge from j to i then we say the edge from i to j is *reciprocated*. (Obviously the edge from j to i is also reciprocated.) Pairs of edges like this are also sometimes called *co-links*, particularly in the context of the World Wide Web [157]. The reciprocity r is defined as the fraction of edges that are reciprocated. Noting that the product of adjacency matrix elements $A_{ij}A_{ji}$ is 1 if and only if there is an edge from i to j and an edge from j to i and is zero otherwise, we can sum over all node pairs i, j to get an expression for the reciprocity:

$$r = \frac{1}{m} \sum_{ij} A_{ij}A_{ji} = \frac{1}{m} \text{Tr } \mathbf{A}^2, \quad (7.32)$$

where m is, as usual, the total number of (directed) edges in the network.

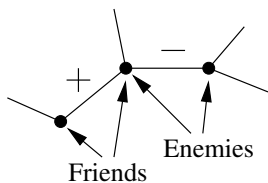
Consider, for example, this small network of four nodes:



A loop of length two in a directed network.

There are seven directed edges in this network and four of them are reciprocated, so the reciprocity is $r = \frac{4}{7} \approx 0.57$. In fact, this is about the same value as is seen on the World Wide Web. There is about a 57% percent chance that if web page A links to web page B then B also links back to A.²⁰ As another example, in a network of who has whom in their email address book it was found that the reciprocity was about $r = 0.23$ [364]. And in a study of friendship networks from a large set of US high schools, reciprocity values were found to range from about 0.3 to 0.5, depending on the school [38].

7.5 SIGNED EDGES AND STRUCTURAL BALANCE



Friends and enemies in an acquaintance network can be denoted by positive and negative edges.

In some social networks, and occasionally in other networks, edges are allowed to be either “positive” or “negative.” For instance, in an acquaintance network we could denote friendship by a positive edge and animosity by a negative edge. One could also consider varying degrees of friendship or animosity—networks with more strongly positive or negative edges in them—but for the moment let’s stick to the simple case where each edge is in just one of two states, positive or negative, like or dislike. Such networks are sometimes called *signed networks* and their edges are called *signed edges*.

It is important to be clear that a negative edge in this context is not the same as the absence of an edge. A negative edge indicates, for example, two people who are acquainted but dislike each other. The absence of an edge represents two people who are not acquainted—whether they would like one another if they ever met is not recorded.

Now consider the possible configurations of three edges in a triangle in a signed network, as depicted in Fig. 7.9. If “+” and “-” represent like and dislike, then we can imagine some of these configurations creating social problems if they were to arise between three people in the real world. Configuration (a) is fine: everyone likes everyone else. Configuration (b) is probably also fine, although the situation is more subtle than (a). Individuals u and v like one another and both dislike w , but the configuration can still be regarded as stable in the sense that u and v can agree over their dislike of w and get along just fine, while w hates both of them. No one is conflicted about their allegiances. Put another way, there is no problem with u and v being friends if one considers that “the enemy of my enemy is my friend.”

²⁰This figure is an unusually high one among directed networks, but there are reasons for it. One is that many of the links between web pages are between pages on the same website, and it is common for such pages to link to each other. If you were to exclude links between pages on the same site the value of the reciprocity would certainly be lower.

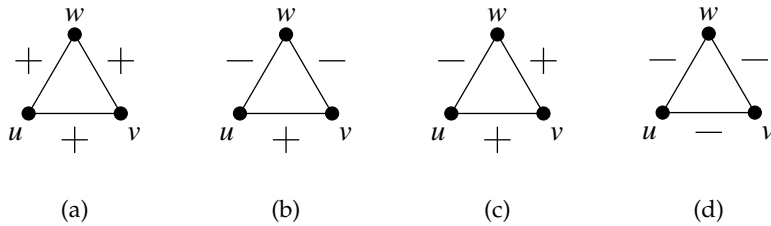


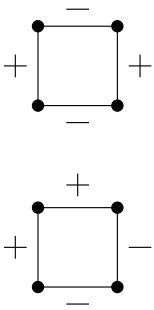
Figure 7.9: Possible triad configurations in a signed network. Configurations (a) and (b) are balanced and hence relatively stable, but configurations (c) and (d) are unbalanced and liable to break apart.

Configuration (c) however could be problematic. Individual u likes individual v and v likes w , but u thinks w is an idiot. This is going to place a strain on the friendship between u and v because u thinks v 's friend is an idiot. Alternatively, from the point of view of v , v has two friends, u and w , and they don't get along, which puts v in an awkward position. In many real-life situations this kind of tension would be resolved by one of the acquaintanceships being broken; i.e., the edge would be removed altogether. Perhaps v would simply stop talking to one of his friends, for instance.

Configuration (d) is somewhat ambiguous. On the one hand, it consists of three people who all dislike each other, so no one is in doubt about where things stand: everyone just hates everyone else. On the other hand, the "enemy of my enemy" rule is broken here. Individuals u and v might like to form an alliance in recognition of their joint dislike of w , but find it difficult to do so because they also dislike each other. This could cause tension—think of the uneasy alliance of the US and Russia against Germany during World War II, for instance. But what one can say definitely is that configuration (d) is often unstable. There may be little reason for the three to stay together when none of them likes the others. Quite probably three enemies such as these would simply sever their connections and go their separate ways.

The feature that distinguishes the two stable configurations in Fig. 7.9 from the unstable ones is that they have an even number of minus signs around the loop.²¹ One can enumerate similar configurations for longer loops, of length four or greater, and again find that loops with even numbers of minus signs

²¹This is similar in spirit to the concept of "frustration" that arises in the physics of magnetic spin systems.



Two stable configurations in loops of length four.

appear stable and those with odd numbers unstable.

This alone would be an observation of only slight interest, were it not for the intriguing fact that this type of stability does appear to extend to real networks. In surveys it is found that the unstable configurations in Fig. 7.9, those with odd numbers of minus signs, occur far less often in real social networks than the stable configurations with even numbers of minus signs [29, 167]. Networks containing only loops with even numbers of minus signs are said to show *structural balance*, or sometimes just *balance*.

An important consequence of structural balance in networks was proved by Harary [229]:

A balanced network can be divided into connected groups of nodes such that all connections between members of the same group are positive and all connections between members of different groups are negative.

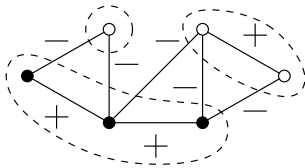


Figure 7.10: A balanced, clusterable network. Every loop in this network contains an even number of minus signs. The dotted lines indicate the division of the network into clusters such that all acquaintances within clusters have positive connections and all acquaintances in different clusters have negative connections.

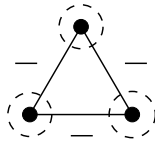
Note that the groups in question can consist of a single node or many nodes, and there may be only one group or there may be very many. Figure 7.10 shows a balanced network and its division into groups. Networks that can be divided into groups like this are said to be *clusterable*. Harary’s balance theorem tells us that all balanced networks are clusterable.

The theorem is straightforward to prove. Imagine coloring the nodes of the network with two colors, such that nodes at the ends of a positive edge are always the same color and nodes at the ends of a negative edge are different colors. It’s not hard to show that this is always possible if all loops in the network have an even number of minus signs—see Exercise 7.13 on page 216. Once the nodes have been colored in this way, we can immediately deduce the identity of the groups that satisfy Harary’s theorem: we simply divide the network into contiguous clusters of nodes that have the same color—see Fig. 7.10. In every such cluster, since all nodes have the same color, they must be joined by positive edges, while at the same time all edges that connect different clusters must be negative, since the clusters have different colors—if they did not have different colors (and were connected by at least one edge) they would be considered the same cluster.²²

²²As an interesting historical note, we observe that while Harary’s theorem is perfectly correct, his interpretation of it was, in this author’s opinion, erroneous. In his 1953 paper [229], he describes the meaning of the theorem in the following words: “A psychological interpretation of Theorem 1 is that a ‘balanced group’ consists of two highly cohesive cliques which dislike each other.” (Harary

The practical importance of Harary’s result rests on the fact that, as mentioned earlier, many real social networks are naturally found to be in a balanced or mostly balanced state. In such cases it would be possible, therefore, for the network to divide into groups such that everyone likes the people they know in their own group and dislikes those in other groups (or nearly so, in the case of an approximately but not perfectly balanced network). Structural balance and clusterability is thus a model for the evolution of cliquishness or insularity, with people tending to stick together in like-minded groups and disdaining everyone outside their immediate community.

We can also ask whether the inverse of Harary’s theorem is true. Is it also the case that a network that is clusterable is necessarily balanced? The answer is no, as this simple counter-example shows:



In this network all three nodes dislike each other, so there are an odd number of minus signs around the loop, but there is no problem dividing the network into three clusters of one node each such that everyone dislikes the members of the other clusters. Thus this network is clusterable but not balanced. Davis [132] proved that a necessary and sufficient condition for clusterability is that the network contain no loops with exactly one minus sign. The network above clearly contains no such loops, and hence is indeed clusterable. The proof of Davis’s result is based on a generalization of Harary’s theorem to the case where we color the nodes with more than two colors—see Exercise 7.14 on page 216.

Real networks are not always perfectly balanced or clusterable, but nonethe-

is using the word “clique” in a non-technical sense here to mean a closed group of people, rather than in the graph theoretical sense of Section 7.2.1.) However, just because it is possible to color the network with two colors as described above does not mean the network forms two groups. Since the nodes of a single color are not necessarily contiguous, there are in general many groups of each color, and it seems unreasonable to describe these groups as forming a “highly cohesive clique” when in fact they have no contact at all. Moreover, it is neither possible nor correct to conclude that the members of two groups of opposite colors dislike each other unless there is at least one edge connecting the two. If two groups of opposite colors never actually have any contact then it might be that they would get along just fine if they met. It is straightforward to show that such an occurrence would lead to an unbalanced network, but Harary’s statement says that the *present* balanced network implies dislike, and this is untrue. Only if the network were to remain balanced upon addition of one or more edges between groups of unlike colors would his conclusion be accurate.

less the basic ideas may still apply. It may not be possible to divide the nodes of a network into groups such that all internal connections are positive, but it is often possible to find a division where most are positive. For example, Axelrod and Bennett [33] studied the pattern of alignments—cordial versus hostile—between 17 European countries on the eve of the Second World War. We can think of this as a network of country nodes connected by positive and negative edges. They looked for divisions of the network into groups such that there were few hostile interactions within groups and found that the best such division, in a sense that they defined, corresponded closely to the actual division of powers during the war, with Germany, Italy, and others on one side and Britain, France, the Soviet Union, and their allies on the other. Only two countries, Poland and Portugal, were incorrectly assigned by the calculation.²³

7.6 SIMILARITY

Another central concept in social network analysis is that of similarity between nodes. In what ways are nodes in a network similar, and how can we quantify that similarity? Which nodes in a given network are most similar to one another? Which node v is most similar to a given node u ? Answers to questions like these can help us tease apart the relationships between nodes in social networks, information networks, and others. It might be useful, for instance, to have a list of web pages that are similar—in some appropriate sense—to another page that we specify. In fact, some web search engines already provide a feature like this: “Click here for pages similar to this one.”

Similarity can be determined in many different ways and most of them have nothing to do with networks. For example, commercial dating and match-making services try to match people with others to whom they are similar by using descriptions of people’s interests, background, likes, and dislikes. In effect, these services are computing similarity measures between people based on personal characteristics. Our focus in this book, however, is on networks, so we will concentrate on the more limited problem of determining similarity between the nodes of a network using the information contained in the network structure.

There are two fundamental approaches to constructing measures of network similarity, called *structural equivalence* and *regular equivalence*. The names are

²³The situation is complicated by the fact that the countries in question had widely varying sizes. The calculation performed by Axelrod and Bennett took the sizes of the countries into account as well as their alignments, with larger countries carrying more weight than smaller ones, so it is a more complex computation than simply counting edges.

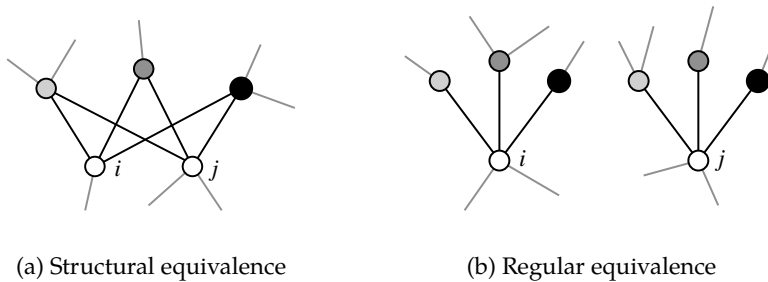


Figure 7.11: Structural equivalence and regular equivalence. (a) Nodes i and j are structurally equivalent if they share many of the same neighbors. (b) Nodes i and j are regularly equivalent if their neighbors are themselves equivalent (indicated here by the different shades of nodes).

rather opaque, but the ideas they represent are simple enough. Two nodes in a network are structurally equivalent if they share many of the same network neighbors. In Fig. 7.11a we show a sketch depicting structural equivalence between two nodes i and j —the two share, in this case, three of the same neighbors, although both also have other neighbors that are not shared.

Regular equivalence is more subtle. Two regularly equivalent nodes do not necessarily share the same neighbors, but they have neighbors *who are themselves similar*. Two history students at different universities, for example, may not have any friends in common, but they can still be similar in the sense that they both know a lot of other history students, history instructors, and so forth. Similarly, two CEOs at different companies may have no colleagues in common, but they are similar in the sense that they have professional ties to their respective CFO, CIO, members of the board, company president, and so forth. Regular equivalence is illustrated in Fig. 7.11b.

7.6.1 MEASURES OF STRUCTURAL EQUIVALENCE

Perhaps the simplest and most obvious measure of structural equivalence is just a count of the number of common neighbors two nodes have. In an undirected network the number n_{ij} of common neighbors of nodes i and j is given by

$$n_{ij} = \sum_k A_{ik}A_{kj}, \quad (7.33)$$

which is just the ij th element of \mathbf{A}^2 .

However, a simple count of common neighbors is not on its own a very good measure of similarity. If two nodes have three common neighbors, is that a lot or a little? It's hard to tell unless we know, for instance, what the degrees of the nodes are, or how many common neighbors other pairs of nodes share. What we need is some sort of normalization that places the similarity value on an easily understood scale. One strategy might be to divide by the total number of nodes in the network n , since this is the maximum number of common neighbors two nodes can have in a simple graph. (Technically, the maximum is actually $n - 2$, but the difference is small when n is large.) However, this unduly penalizes nodes with low degree: if a node has degree three then it can have at most three neighbors in common with another node, but the two nodes would still receive a small similarity value in the common case where n is large. A better measure would allow for the varying degrees of nodes. Such a measure is the *cosine similarity*, sometimes also called *Salton's cosine*.

In geometry, the inner or dot product of two vectors \mathbf{x} and \mathbf{y} is given by $\mathbf{x} \cdot \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos \theta$, where $|\mathbf{x}|$ and $|\mathbf{y}|$ are the magnitudes of the two vectors and θ is the angle between them. Rearranging, we can write the cosine of the angle as

$$\cos \theta = \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}. \quad (7.34)$$

Salton [422] proposed that we regard the i th and j th rows (or columns) of the adjacency matrix as two vectors and use the cosine of the angle between them as our similarity measure. Noting that the dot product of two rows is simply $\sum_k A_{ik} A_{kj}$ for an undirected network, this gives us a similarity

$$\sigma_{ij} = \cos \theta = \frac{\sum_k A_{ik} A_{kj}}{\sqrt{\sum_k A_{ik}^2} \sqrt{\sum_k A_{jk}^2}}. \quad (7.35)$$

Assuming our network is an unweighted simple graph, the elements of the adjacency matrix take only the values 0 and 1, so that $A_{ij}^2 = A_{ij}$ for all i, j . Then $\sum_k A_{ik}^2 = \sum_k A_{ik} = k_i$, where k_i is the degree of node i (see Eq. (6.12)). Thus

$$\sigma_{ij} = \frac{\sum_k A_{ik} A_{kj}}{\sqrt{k_i} \sqrt{k_j}} = \frac{n_{ij}}{\sqrt{k_i k_j}}. \quad (7.36)$$

In other words, the cosine similarity of i and j is the number of common neighbors of the two nodes divided by the geometric mean of their degrees. For the nodes i and j depicted in Fig. 7.11a, for instance, the cosine similarity would be

$$\sigma_{ij} = \frac{3}{\sqrt{4 \times 5}} = 0.671 \dots \quad (7.37)$$

Note that the cosine similarity is technically undefined if one or both of the nodes has degree zero, but by convention we normally say in that case that $\sigma_{ij} = 0$.

The cosine similarity provides a natural scale for quantifying similarity. Its value always lies in the range from zero to one. A cosine similarity of one indicates that two nodes have exactly the same neighbors. A cosine similarity of zero indicates that they have none of the same neighbors. Note that the cosine similarity can never be negative, being a sum of positive terms, even though cosines of angles in general can be negative.

Cosine similarity is the most widely used measure of similarity in networks but not the only one. Another common measure is the *Jaccard coefficient*, which similarly normalizes n_{ij} to run between zero and one, but does so in a slightly different fashion. The Jaccard coefficient for nodes i and j is defined to be the number of common neighbors n_{ij} divided by the total number of distinct neighbors of both nodes. That is, the normalizing factor is the total number of neighbors of both, but counting common neighbors only once, not twice, for a total of $k_i + k_j - n_{ij}$ neighbors. Thus the Jaccard coefficient is

$$J_{ij} = \frac{n_{ij}}{k_i + k_j - n_{ij}}. \quad (7.38)$$

When i and j have no neighbors in common, so that $n_{ij} = 0$, this gives $J_{ij} = 0$. When they have all of their neighbors in common and $k_i = k_j = n_{ij}$ it gives $J_{ij} = 1$. In all other situations it gives a value somewhere in between.

A number of other similarity measures find occasional use as well, including the Pearson correlation coefficient between rows of the adjacency matrix:

$$r_{ij} = \frac{\sum_k (A_{ik} - \langle A_i \rangle)(A_{jk} - \langle A_j \rangle)}{\sqrt{\sum_k (A_{ik} - \langle A_i \rangle)^2} \sqrt{\sum_k (A_{jk} - \langle A_j \rangle)^2}}, \quad (7.39)$$

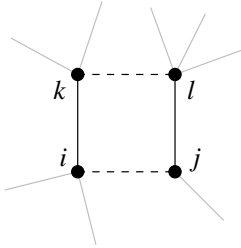
where $\langle A_i \rangle$ is the average of the i th row, and the *Hamming distance* (sometimes also called *Euclidean distance*), which is the number of neighbors two nodes *don't* share in common (i.e., the number that are neighbors of one node but not the other):

$$h_{ij} = \sum_k (A_{ik} - A_{jk})^2. \quad (7.40)$$

Technically, Hamming distance is really a *dissimilarity* measure, since it is larger for nodes with fewer common neighbors. For further discussion of these measures see, for instance, Wasserman and Faust [462].

7.6.2 MEASURES OF REGULAR EQUIVALENCE

The other type of similarity considered in social network analysis is regular equivalence. As described in Section 7.6.1, regularly equivalent nodes are nodes that, while they do not necessarily share neighbors, have neighbors who are themselves similar—see Fig. 7.11b again.



Nodes i and j are considered similar (dashed line) if they have respective neighbors k and l that are themselves similar.

Quantitative measures of regular equivalence are less well developed than measures of structural equivalence, but a number of measures have been proposed in recent years that appear to work reasonably well. The basic idea [65,248,296] is to define a similarity score σ_{ij} such that i and j have high similarity if they have neighbors k and l that themselves have high similarity. For an undirected network we can write this as

$$\sigma_{ij} = \alpha \sum_{kl} A_{ik}A_{jl}\sigma_{kl} , \tag{7.41}$$

or in matrix terms $\sigma = \alpha \mathbf{A}\sigma\mathbf{A}$, where α is a constant. Although it may not be immediately apparent, this expression is a type of eigenvector equation, where the entire matrix σ of similarities is the eigenvector. The parameter α is the eigenvalue (or more correctly, its inverse) and, as with the eigenvector centrality of Section 7.1.2, we are normally interested in the leading eigenvector.

This formula however has some problems. First, it doesn't necessarily give a high value for the "self-similarity" σ_{ii} of a node to itself, which is counter-intuitive since presumably all nodes are similar to themselves. As a consequence of this, Eq. (7.41) also doesn't necessarily give a high similarity score to node pairs that have a lot of common neighbors, which in the light of our discussion of structural equivalence in the preceding section we might feel it should. If we had high self-similarity scores for all nodes, then Eq. (7.41) would automatically give high similarity also to nodes with many common neighbors, because for such nodes the sum on the right-hand side would have large contributions from terms of the form $A_{ik}A_{jk}\sigma_{kk}$.

We can address these problems by introducing an extra diagonal term in the similarity thus:

$$\sigma_{ij} = \alpha \sum_{kl} A_{ik}A_{jl}\sigma_{kl} + \delta_{ij} , \tag{7.42}$$

or in matrix notation

$$\sigma = \alpha \mathbf{A}\sigma\mathbf{A} + \mathbf{I} , \tag{7.43}$$

which gives an extra boost to the similarity score of a node with itself.

This still has some problems though. Suppose we evaluate Eq. (7.43) by repeated iteration, taking a starting value, for example, of $\sigma^{(0)} = 0$ and using

it to compute $\sigma^{(1)} = \alpha \mathbf{A} \sigma^{(0)} \mathbf{A} + \mathbf{I}$, and then repeating the process many times until σ converges. On the first few iterations we get the results

$$\sigma^{(1)} = \mathbf{I}, \quad (7.44a)$$

$$\sigma^{(2)} = \alpha \mathbf{A}^2 + \mathbf{I}, \quad (7.44b)$$

$$\sigma^{(3)} = \alpha^2 \mathbf{A}^4 + \alpha \mathbf{A}^2 + \mathbf{I}. \quad (7.44c)$$

The pattern is clear: in the limit of many iterations, we will get a sum over even powers of the adjacency matrix. However, as discussed in Section 6.11, the elements of the r th power of the adjacency matrix count paths of length r between nodes, and hence this measure of similarity is a weighted sum over paths of even length between pairs of nodes.

But why should we consider only paths of even length? Why not consider paths of all lengths? These questions lead us to a better definition of regular equivalence as follows:²⁴ nodes i and j are similar if i has a neighbor k that is itself similar to j . Again we assume that nodes are similar to themselves, which we can represent with a diagonal δ_{ij} , and our similarity measure then looks like

$$\sigma_{ij} = \alpha \sum_k A_{ik} \sigma_{kj} + \delta_{ij}, \quad (7.45)$$

or

$$\sigma = \alpha \mathbf{A} \sigma + \mathbf{I}, \quad (7.46)$$

in matrix notation. Evaluating this expression by again iterating from a starting value of $\sigma^{(0)} = 0$, we get

$$\sigma^{(1)} = \mathbf{I}, \quad (7.47a)$$

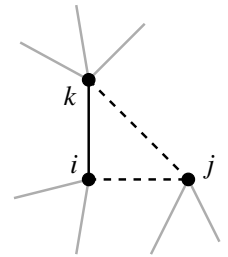
$$\sigma^{(2)} = \alpha \mathbf{A} + \mathbf{I}, \quad (7.47b)$$

$$\sigma^{(3)} = \alpha^2 \mathbf{A}^2 + \alpha \mathbf{A} + \mathbf{I}. \quad (7.47c)$$

In the limit of a large number of iterations this gives

$$\sigma = \sum_{m=0}^{\infty} (\alpha \mathbf{A})^m = (\mathbf{I} - \alpha \mathbf{A})^{-1}, \quad (7.48)$$

which we could also have deduced directly by rearranging Eq. (7.46). Now our similarity measure includes counts of paths of all lengths, not just even-length paths. In fact, we can see now that this similarity measure could be defined



In the modified definition of regular equivalence node i is considered similar to node j (dashed line) if it has a neighbor k that is itself similar to j .

²⁴This definition is not obviously symmetric with respect to i and j but, as we will see, does in fact give rise to a symmetric expression.

in a completely different way, as a weighted count of all the paths between the nodes i and j with paths of length r getting weight α^r . As long as $\alpha < 1$, longer paths will get less weight than shorter ones, which seems sensible. In effect we are saying that nodes are similar if they are connected either by a few short paths or by very many long ones.

Equation (7.48) is reminiscent of the formula for the Katz centrality, Eq. (7.7). We could call Eq. (7.48) the “Katz similarity” perhaps, although Katz himself never discussed it. It has the nice property that the Katz centrality of a node is equal to the sum of the Katz similarities of that node to all others, so that nodes that are similar to many others would get high centrality, a concept that certainly makes intuitive sense. As with the Katz centrality, the value of the parameter α is undetermined—we are free to choose it as we see fit—but it must satisfy $\alpha < 1/\kappa_1$ if the sum in Eq. (7.48) is to converge, where κ_1 is the largest eigenvalue of the adjacency matrix. (See the discussion in Section 7.1.3.)

In a sense, this regular equivalence measure can be seen as a generalization of structural equivalence measures such as the cosine similarity and Jaccard coefficient of the preceding section. Those measures were based on a count of the number of common neighbors of a pair of nodes, but the number of common neighbors is also equal to the number of paths of length two. Our “Katz similarity” measure simply extends this approach to paths of all lengths.

Some variations of the Katz similarity are possible. As defined it tends to give high similarity to nodes that have high degree, because high-degree nodes have more terms in the sum in Eq. (7.45). In some cases this might be desirable: maybe the person with many friends *should* be considered more similar to others than the person with few. However, in other cases it gives an unwanted bias in favor of high-degree nodes. Who is to say that two hermits are not “similar” in an interesting sense? If we wish, we can remove the bias in favor of high degree by dividing by node degree thus:

$$\sigma_{ij} = \frac{1}{k_i} \left[\alpha \sum_k A_{ik} \sigma_{kj} + \delta_{ij} \right], \quad (7.49)$$

or in matrix notation $\sigma = \mathbf{D}^{-1}(\alpha \mathbf{A} \sigma + \mathbf{I})$, where, as previously, \mathbf{D} is the diagonal matrix with elements $D_{ii} = k_i$. This expression can be rearranged to read

$$\sigma = (\mathbf{D} - \alpha \mathbf{A})^{-1}. \quad (7.50)$$

In much the same way that the Katz similarity is related to Katz centrality, this similarity measure is related to PageRank, though the correspondence is not perfect: the sum of centralities σ_{ij} calculated from Eq. (7.50) over all nodes j does not give the PageRank of node i ; it gives the PageRank divided by k_i —see Exercise 7.15.

Another variant allows for cases where the last term in Eqs. (7.45) or (7.49) is not simply diagonal, but includes off-diagonal elements too. This would allow us, for example, to specify explicitly that particular pairs of nodes are similar, based on some other (probably non-network) information that we have at our disposal. Going back to the example of CEOs at companies that we gave at the beginning of Section 7.6, we could, for instance, specify that the CFOs and CIOs and so forth at different companies are similar, and then our similarity measure would, we hope, correctly deduce from the network structure that the CEOs are similar also. This kind of approach is particularly useful in the case of networks that consist of more than one component, so that some pairs of nodes are not connected at all. If, for example, we have two separate components representing two different companies, then there will be no paths of any length between individuals in different companies, and hence a measure like (7.45) or (7.49) will never assign a non-zero similarity to such individuals. But if we explicitly assert some similarities between members of the different companies, our measure will then be able to build on that information to deduce similarities between other members.²⁵

7.7 HOMOPHILY AND ASSORTATIVE MIXING

Consider Fig. 7.12, which shows a friendship network of students at an American high school, determined from a questionnaire of the type discussed in Section 4.2.²⁶ One clear feature that emerges from the figure is the division of the network into two groups. It turns out in this case that the division is principally along lines of race: the different shades of the nodes in the picture correspond to students of different race as denoted in the legend, and reveal that the school is divided between a group composed principally of black students and a group composed principally of white.

This is not news to sociologists, who have long observed and discussed such divisions [338]. Nor is the effect specific to race. People are found to form friendships, acquaintances, business relations, and many other types of

²⁵The idea of generalizing from a few given inputs to a whole system in this way is common in fields like machine learning and information retrieval. For instance, there is a considerable literature on how to classify objects such as text documents into topics or groups by combining textual or other clues with a small set of initial group assignments. Such problems fall within the general area known as *semi-supervised learning* [100].

²⁶The study used a “name generator”—students were asked to list the names of others they considered to be their friends. This results in a directed network, but we have neglected the edge directions in the figure. In our representation there is an undirected edge between two people if either of them considers the other to be a friend (or both do).

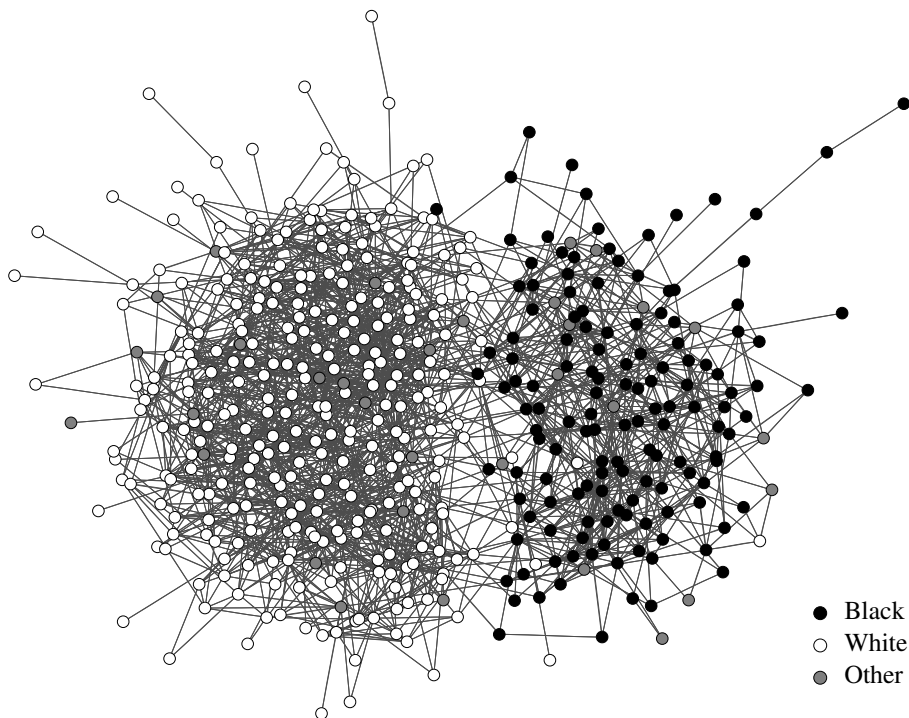


Figure 7.12: A friendship network at a US high school. The nodes in this network represent 470 students at a US high school (ages 14 to 18 years). The nodes are color coded by race as indicated in the key. Data are from the National Longitudinal Study of Adolescent Health [52,451].

ties based on all sorts of characteristics, including age, nationality, language, income, educational level, and others. Almost any social parameter you can imagine plays into people's selection of their friends. People have, it appears, a strong tendency to associate with others whom they perceive as being similar to themselves in some way. This tendency is called *homophily* or *assortative mixing*.

More rarely, one also encounters *disassortative mixing*, the tendency for people to associate with others who are *unlike* them. Probably the most widespread and familiar example of disassortative mixing is mixing by gender in sexual contact networks. The majority of sexual partnerships are between individuals of opposite sex, so they represent connections between people who differ in their gender. Of course, same-sex partnerships do also occur, but they are a smaller fraction of the ties in the network.

Assortative (or disassortative) mixing is also seen in some non-social net-

works. In citation networks, for instance, papers tend to cite other papers in the same field more than they do papers in different fields. Similarly, web pages written in a particular language tend to link to others in the same language.

In this section we look at how assortative mixing can be quantified. Assortative mixing by unordered characteristics such as race, gender, or nationality is fundamentally different from mixing by ordered characteristics like age or income, so we treat the two cases separately.

7.7.1 ASSORTATIVE MIXING BY UNORDERED CHARACTERISTICS

Suppose we have a network in which the nodes are classified according to some characteristic that has a finite set of possible values. The values are descriptive only and don't fall in any particular order. For instance, the nodes of the network could represent people and be classified according to nationality, race, or gender. Or they could be web pages classified by what language they are written in, or biological species classified by habitat, or any of many other possibilities.

The network is assortative if a significant fraction of the edges in the network run between nodes of the same type. One simple way to quantify assortativity would be just to record this fraction, but this is not a very good measure because, for instance, it is 1 if all nodes belong to the same single type. This is a trivial sort of assortativity: all friends of a human being, for example, are also human beings,²⁷ but this is not really a useful statement. What we would like instead is a measure that is large in non-trivial cases but small in trivial ones.

A better measure turns out to be the following. We find the fraction of edges that run between nodes of the same type, and then we subtract from that figure the fraction of such edges we would *expect* to find if edges were positioned at random without regard for node type. For the trivial case in which all nodes are of a single type, for instance, 100% of edges run between nodes of the same type, but this is also the expected figure if edges were placed at random, since there is nowhere else for the edges to fall. The difference of the two numbers is then zero, telling us that there is no non-trivial assortativity in this case. Only when the fraction of edges between nodes of the same type is significantly greater than we would get if the edges were randomly placed will our measure give a large score. Thus, this measure is in a sense quantifying the level of non-randomness in the placement of edges in the network.

In mathematical terms this measure can be written as follows. Let us denote by g_i the group, class, or type of node i , which is an integer $g_i = 1 \dots N$, with

²⁷Ignoring, for the purposes of argument, dogs, cats, imaginary friends, and so forth.

N being the total number of groups. Then the total number of edges that run between nodes of the same type—the number of edges within groups—is

$$\sum_{\text{edges } (i,j)} \delta_{g_i g_j} = \frac{1}{2} \sum_{ij} A_{ij} \delta_{g_i g_j}, \quad (7.51)$$

where δ_{ij} is the Kronecker delta and the factor of $\frac{1}{2}$ compensates for the fact that every node pair i, j is counted twice in the second sum.

Calculating the expected number of edges between nodes if edges are placed at random takes a little more work. Consider a particular node i with degree k_i , and consider a particular edge attached to that node. There are by definition $2m$ ends of edges in the entire network, where m is as usual the total number of edges, and the chances that the other end of our particular edge is one of the k_j ends attached to node j is thus $k_j/2m$ if connections are made purely at random.^{28,29} Counting all k_i edges attached to i , the total expected number of edges between nodes i and j is then $k_i k_j / 2m$, and the expected number of edges between all pairs of nodes of the same type is

$$\frac{1}{2} \sum_{ij} \frac{k_i k_j}{2m} \delta_{g_i g_j}, \quad (7.52)$$

where the factor of $\frac{1}{2}$, as before, compensates for the double counting of node pairs. Subtracting (7.52) from (7.51) then gives us the difference between the actual and expected number of edges in the network that join nodes of the same type:

$$\frac{1}{2} \sum_{ij} A_{ij} \delta_{g_i g_j} - \frac{1}{2} \sum_{ij} \frac{k_i k_j}{2m} \delta_{g_i g_j} = \frac{1}{2} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{g_i g_j}. \quad (7.53)$$

Conventionally, one calculates not the number of such edges but the fraction, which is given by this same expression divided by the total number m of edges. The resulting quantity is called the *modularity* [352, 366], usually denoted Q :

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{g_i g_j}. \quad (7.54)$$

²⁸The exact expression is actually $k_j/(2m-1)$, since we know that one end of one of the edges is definitely attached to node i . In any but the smallest of networks, however, m is a large number and $k_j/2m$ is a good approximation.

²⁹Technically, we are making connections at random while preserving the node degrees. We could in principle ignore node degrees and make connections truly at random, but in practice this is found to give poor results.

Modularity is a measure of the extent to which like is connected to like in a network. It is strictly less than 1 and takes positive values if there are more edges between nodes of the same type than we would expect by random chance. It can also take negative values if there are fewer such edges than we would expect by chance.

For Fig. 7.12, for instance, where the nodes are of three types according to ethnicity—"black," "white," and "other"—we find a positive modularity value of $Q = 0.305$, indicating assortative mixing by race in this particular network. Negative values of the modularity indicate disassortative mixing. We might see a negative modularity, for example, in a network of sexual partnerships where most partnerships are between individuals of opposite sex.

An alternative form for the modularity, which is useful for certain kinds of calculations, can be derived in terms of the quantities

$$e_r = \frac{1}{2m} \sum_{ij} A_{ij} \delta_{g_i,r} \delta_{g_j,r}, \quad (7.55)$$

which is the fraction of edges that join nodes of type r , and

$$a_r = \frac{1}{2m} \sum_i k_i \delta_{g_i,r}, \quad (7.56)$$

which is the fraction of ends of edges attached to nodes of type r . Noting that

$$\delta_{g_i g_j} = \sum_r \delta_{g_i,r} \delta_{g_j,r}, \quad (7.57)$$

we have, from Eq. (7.54),

$$\begin{aligned} Q &= \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \sum_r \delta_{g_i,r} \delta_{g_j,r} \\ &= \sum_r \left[\frac{1}{2m} \sum_{ij} A_{ij} \delta_{g_i,r} \delta_{g_j,r} - \frac{1}{2m} \sum_i k_i \delta_{g_i,r} \frac{1}{2m} \sum_j k_j \delta_{g_j,r} \right] \\ &= \sum_r (e_r - a_r^2). \end{aligned} \quad (7.58)$$

This form can be useful, for instance, when we have network data in the form of a list of edges and the types of the nodes at their ends, but no explicit data on node degrees. In such a case e_r and a_r are relatively easy to calculate, while Eq. (7.54) is quite awkward. See Exercise 7.16 on page 216 for some examples.

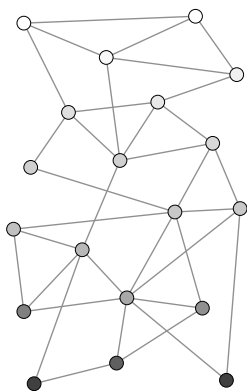
7.7.2 ASSORTATIVE MIXING BY ORDERED CHARACTERISTICS

We can also have assortative mixing in a network according to characteristics like age or income, whose values come in a particular order, so that it is possible to say not only when two nodes are exactly the same according to the characteristic but also when they are approximately the same. For instance, while two people can certainly be of exactly the same age—born on the same day even—they can also be approximately the same age—born within a couple of years of one another, say—and people could (and in fact often do) choose whom they associate with on the basis of such approximate ages. There is no equivalent approximate similarity for the unordered characteristics of the previous section: there is no sense in which people from France and Germany, say, are more nearly of the same nationality than people from France and Spain.³⁰

If network nodes with similar values of a scalar characteristic tend to be connected together more often than those with different values, then the network is considered assortatively mixed according to that characteristic. If, for example, people are friends with others around the same age as them, then the network is assortatively mixed by age. Sometimes you also hear it said that such a network is *stratified* by age, which means the same thing—one can think of age as a one-dimensional scale or axis, with individuals of different ages forming connected “strata” within the network.

Consider Fig. 7.13, which shows friendship data for the same set of US high school students as Fig. 7.12 but now as a function of age. Each dot in the figure corresponds to one pair of friends and the position of the dot along the two axes gives the ages of the friends, with ages measured by school grades.³¹ As the figure shows, there is substantial assortative mixing by age among the students: many dots lie within the boxes close to the diagonal line that represent friendships between students in the same grade. There is also, in this case, a notable tendency for students to have friends of a wider range of ages as their age increases so there is a lower density of points in the top right box than in the bottom left one.

One could make a crude measure of assortative mixing by scalar characteristics by adapting the ideas of the previous section. One could group the nodes into bins according to the characteristic of interest (say age) and then treat the



A sketch of a stratified network in which most connections run between nodes at or near the same “level” in the network, with level along the vertical axis in this case and denoted by the shades of the nodes.

³⁰One could in principle make up some measure of national differences based say on geographic distance. But if the question we are asked is “Are these two people of the same nationality?” then under normal circumstances the only answers are “yes” and “no.” There is nothing in between.

³¹In the US school system there are 12 numbered grades of one year each and children normally enter the first grade when they are six years old. Thus the ninth grade, for example, corresponds to students of age 14 and 15.

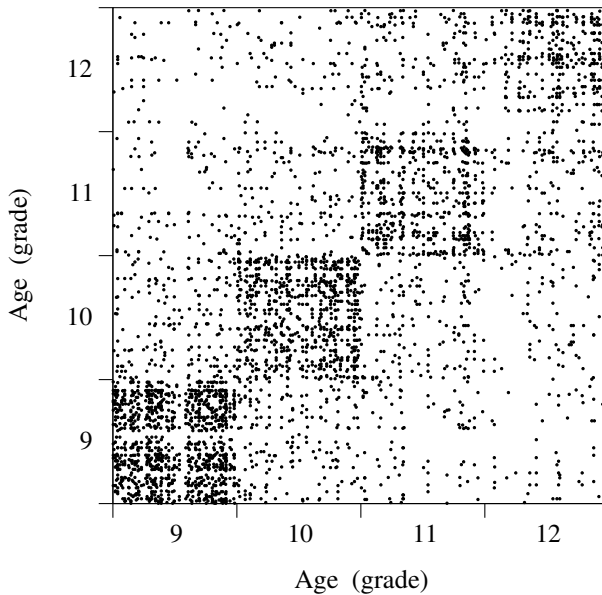


Figure 7.13: Ages of pairs of friends in high school. In this scatter plot each dot corresponds to one of the edges in Fig. 7.12 and its position along the horizontal and vertical axes gives the ages of the two individuals at either end of that edge. The ages are measured in terms of the grades of the students, which run from 9 to 12. In fact, grades in the US school system don't correspond precisely to age since students can start or end their high school careers early or late, and can repeat grades. (Each student is positioned at random within the interval representing their grade, so as to spread the points out on the plot. Note also that each friendship appears twice, once above and once below the diagonal.)

bins as separate types of nodes in the sense of Section 7.7.1. For instance, we might group people by age in ranges of one year or ten years. This, however, somewhat misses the point of scalar characteristics, since it considers nodes falling in the same bin to be of identical types when they may be only approximately so, and nodes falling in different bins to be entirely different when in fact they may be quite similar.

A better approach is to use a covariance measure as follows. Let x_i be the value for node i of the scalar quantity that we are interested in (age, income, etc.). Then consider the pairs of values x_i, x_j for the nodes i, j at the ends of each edge in the network and let us calculate their covariance over all edges as

follows. We define the mean μ of the value of x_i at the end of an edge thus:

$$\mu = \frac{\sum_{ij} A_{ij} x_i}{\sum_{ij} A_{ij}} = \frac{1}{2m} \sum_i k_i x_i, \quad (7.59)$$

where we have made use of Eqs. (6.12) and (6.13). Note that this is not simply the mean value of x_i averaged over all nodes. It is an average over edges, and since a node with degree k_i lies at the ends of k_i edges that node appears k_i times in the average (hence the factor of k_i in the sum).

Then the covariance of x_i and x_j over edges is

$$\begin{aligned} \text{cov}(x_i, x_j) &= \frac{\sum_{ij} A_{ij} (x_i - \mu)(x_j - \mu)}{\sum_{ij} A_{ij}} \\ &= \frac{1}{2m} \sum_{ij} A_{ij} (x_i x_j - \mu x_i - \mu x_j + \mu^2) \\ &= \frac{1}{2m} \sum_{ij} A_{ij} x_i x_j - \mu^2 \\ &= \frac{1}{2m} \sum_{ij} A_{ij} x_i x_j - \frac{1}{(2m)^2} \sum_{ij} k_i k_j x_i x_j \\ &= \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) x_i x_j, \end{aligned} \quad (7.60)$$

where we have made use of Eqs. (6.13) and (7.59). Note the strong similarity between this expression and Eq. (7.54) for the modularity—only the delta function $\delta_{g_i g_j}$ in (7.54) has changed, being replaced by $x_i x_j$.

The covariance will be positive if, on balance, values x_i, x_j at either end of an edge tend to be both large or both small, and negative if they tend to vary in opposite directions. In other words, the covariance will be positive when we have assortative mixing and negative for disassortative mixing.

It is sometimes convenient to normalize the covariance so that it takes the value 1 in a network with perfect assortative mixing—one in which all edges fall between nodes with precisely equal values of x_i (although in most cases such an occurrence would be extremely unlikely). Setting $x_i = x_j$ in the first term of the sum in Eq. (7.60) tells us that the value for a perfectly mixed network would be

$$\frac{1}{2m} \sum_{ij} \left(A_{ij} x_i^2 - \frac{k_i k_j}{2m} x_i x_j \right) = \frac{1}{2m} \sum_{ij} \left(k_i \delta_{ij} - \frac{k_i k_j}{2m} \right) x_i x_j, \quad (7.61)$$

and the normalized measure, sometimes called an *assortativity coefficient*, is the ratio of the two:

$$r = \frac{\sum_{ij}(A_{ij} - k_i k_j / 2m)x_i x_j}{\sum_{ij}(k_i \delta_{ij} - k_i k_j / 2m)x_i x_j}. \quad (7.62)$$

Although it may not be immediately obvious, this is in fact an example of a (Pearson) correlation coefficient—the standard statistical measure of correlation for scalar data—having a covariance in its numerator and a variance in the denominator. The correlation coefficient varies in value between a maximum of 1 for a perfectly assortative network and a minimum of -1 for a perfectly disassortative one. A value of zero implies that the values of x_i at the ends of edges are uncorrelated.³² This normalized correlation coefficient is probably the most widely used measure of assortativity by scalar characteristics.

For the data of Fig. 7.13 the correlation coefficient is found to take a value of $r = 0.616$, indicating that the student friendship network has significant assortative mixing by age—students tend to be friends with others who have ages close to theirs.

It would be possible in principle also to have assortative (or disassortative) mixing according to vector characteristics, with nodes whose vectors have similar values, as measured by some appropriate metric, being more (or less) likely to be connected by an edge. One example of such mixing is the formation of friendships between individuals according to their geographic locations, location being specified by a two-dimensional vector of, for example, latitude/longitude coordinates. It is certainly the case that in general people tend to be friends with others who live geographically close to them, so one would expect mixing of this type to be assortative. Formal treatments of vector assortative mixing, however, have not been much pursued in the networks literature so far.

7.7.3 ASSORTATIVE MIXING BY DEGREE

A special case of assortative mixing according to a scalar quantity, and one of particular interest, is that of mixing by degree. In a network that shows

³²There could be non-linear correlations in such a network and we could still have $r = 0$; the correlation coefficient detects only linear correlations. For instance, we could have nodes with high and low values of x_i connected predominantly to nodes with intermediate values. This is neither assortative nor disassortative by the conventional definition and would give a small value of r , but might nonetheless be of interest. Such non-linear correlations could be discovered by examining a plot such as Fig. 7.13 or by using alternative measures of correlation such as information theoretic measures. Thus, it is perhaps wise not to rely solely on the value of r in investigating assortative mixing.

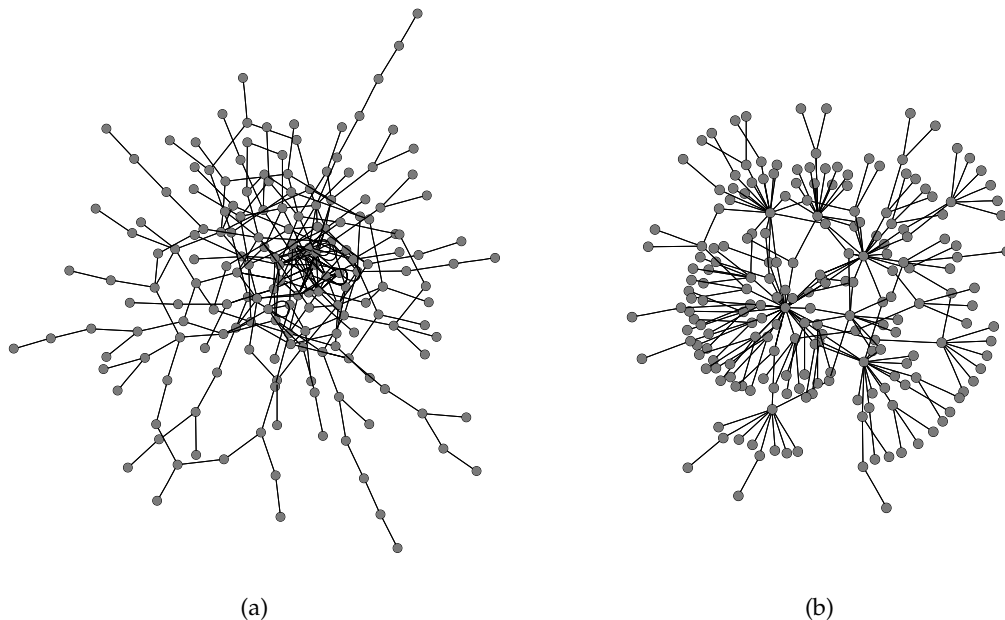


Figure 7.14: Assortative and disassortative networks. These two small networks were computer generated, to illustrate the phenomenon of assortativity by degree. (a) A network that is assortative by degree, displaying the characteristic dense core of high-degree nodes surrounded by a periphery of lower-degree ones. (b) A disassortative network, displaying the star-like structures characteristic of this case. Figure from Newman and Girvan [365]. Copyright 2003 Springer-Verlag Berlin Heidelberg. Reproduced with permission of Springer Nature.

assortative mixing by degree, the high-degree nodes will be preferentially connected to other high-degree nodes, and the low to low. In a social network, for example, we have assortative mixing by degree if the gregarious people are friends with other gregarious people and the hermits with other hermits. Conversely, we could have disassortative mixing by degree, which would mean that the gregarious people were hanging out with the hermits and vice versa.

The reason this case is particularly interesting is because, unlike age or income, degree is itself a property of the network structure. Having one structural property (the degrees) dictate another (the positions of the edges) gives rise to some interesting features in networks. In particular, in an assortative network, where the high-degree nodes tend to stick together, one expects to get a clump or *core* of such high-degree nodes in the network surrounded by a less dense *periphery* of nodes with lower degree. This core–periphery structure is a common feature of many networks, particularly social networks, which are often

We encountered core–periphery structure previously in our discussion of k -cores in Section 7.2.2 and it is discussed further in Section 14.7.3.

found to be assortatively mixed by degree [237, 350, 414]. Figure 7.14a shows a small assortatively mixed network in which the core–periphery structure is clearly visible.

On the other hand, if a network is disassortatively mixed by degree then high-degree nodes tend to be connected to low-degree ones, creating star-like features in the network that are often readily visible. Figure 7.14b shows an example of a small disassortative network. Disassortative networks do not usually have a core–periphery split but are instead more uniform.

Assortative mixing by degree can be measured in the same way as mixing according to any other scalar quantity. We define a covariance of the type described by Eq. (7.60), but with x_i now equal to the degree k_i :

$$\text{cov}(k_i, k_j) = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) k_i k_j, \quad (7.63)$$

or if we wish we can normalize by the maximum value of the covariance to get a correlation coefficient or assortativity coefficient:

$$r = \frac{\sum_{ij} (A_{ij} - k_i k_j / 2m) k_i k_j}{\sum_{ij} (k_i \delta_{ij} - k_i k_j / 2m) k_i k_j}. \quad (7.64)$$

We will see a number of examples of the application of this formula in Section 10.7.

One further thing to note is that evaluating Eq. (7.63) or Eq. (7.64) requires only the structure of the network and no other information, unlike the corresponding calculations for other types of assortative mixing. Once we know the adjacency matrix we also know the degrees of all the nodes and hence we can calculate r . Perhaps for this reason mixing by degree is one of the most frequently studied types of assortative mixing.

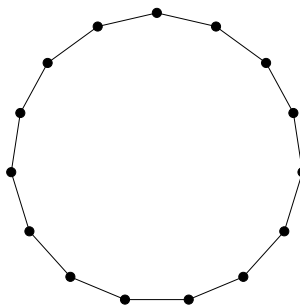
EXERCISES

7.1 Consider a connected k -regular undirected network (i.e., a network in which every node has degree k and there is only one component).

- a) Show that the uniform vector $\mathbf{1} = (1, 1, 1, \dots)$ is an eigenvector of the adjacency matrix with eigenvalue k . In a connected network there is only one eigenvector with all elements positive and hence the eigenvector $\mathbf{1}$ gives, by definition, the eigenvector centrality of our k -regular network and the centralities are the same for every vertex.

- b) Find the Katz centralities of all nodes in the network as a function of k .
- c) You should find that, like the eigenvector centralities, the Katz centralities of all nodes are the same. Name a centrality measure that could give different centrality values for different nodes in a regular network.

7.2 A network consists of n nodes in a ring, where n is odd:



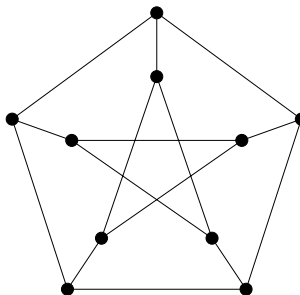
All the nodes have the same closeness centrality. What is it, as a function of n ?

7.3 As we saw in Section 7.1.3, the Katz centrality in vector form satisfies the equation $\mathbf{x} = \alpha \mathbf{A} \mathbf{x} + \mathbf{1}$ (which is Eq. (7.6) with the conventional choice $\beta = 1$).

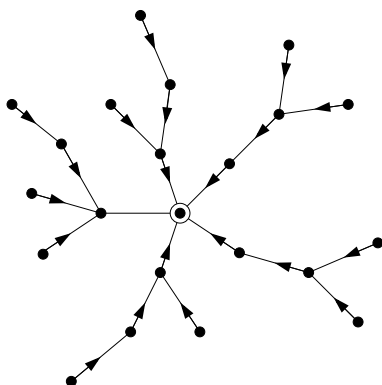
- a) Show that the Katz centrality can also be written in series form as $\mathbf{x} = \mathbf{1} + \alpha \mathbf{A} \mathbf{1} + \alpha^2 \mathbf{A}^2 \mathbf{1} + \dots$
- b) Hence, argue that in the limit where α is small but non-zero, the Katz centrality is essentially equivalent to degree centrality.
- c) Conversely, in the limit $\alpha \rightarrow 1/\kappa_1$, where κ_1 is the largest (most positive) eigenvalue of the adjacency matrix, argue that \mathbf{x} becomes proportional to the leading eigenvector, which is simply the eigenvector centrality.

Thus, the Katz centrality can be thought of as a one-parameter family of centralities, parametrized by $\alpha \in [0, 1/\kappa_1]$, which includes the degree centrality and the eigenvector centrality at the two limits of the range and interpolates between them everywhere in between.

7.4 Calculate the closeness centrality of each of the nodes in this network:

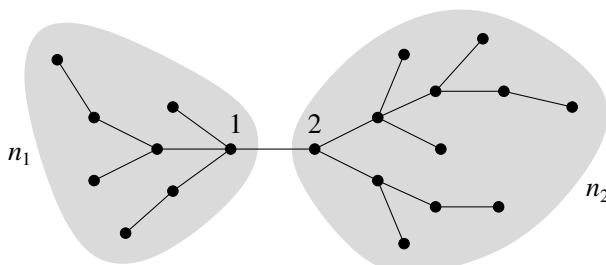


7.5 Suppose a directed network takes the form of a tree with all edges pointing inward towards a central node:



What is the PageRank centrality of the central node in terms of the single parameter α appearing in the definition of PageRank and the distances d_i from each node i to the central node?

7.6 Consider an undirected tree of n nodes. A particular edge in the tree joins nodes 1 and 2 and divides the tree into two disjoint regions of n_1 and n_2 nodes as sketched here:



Show that the closeness centralities C_1 and C_2 of the two nodes, defined according to Eq. (7.21), are related by

$$\frac{1}{C_1} + \frac{n_1}{n} = \frac{1}{C_2} + \frac{n_2}{n}.$$

7.7 Consider an undirected tree of n nodes. Suppose that a particular node in the tree has degree k , so that its removal would divide the tree into k disjoint regions, and suppose that the sizes of those regions are $n_1 \dots n_k$.

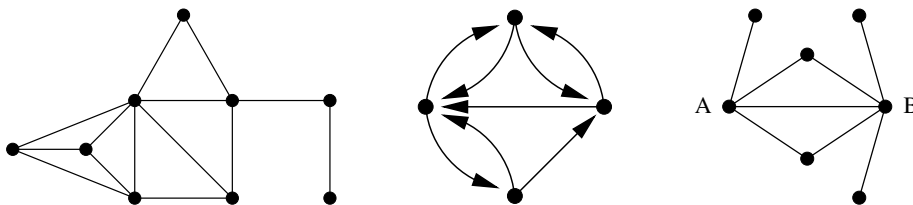
a) Show that the unnormalized betweenness centrality x of this node, as defined in Eq. (7.24), is

$$x = n^2 - \sum_{m=1}^k n_m^2.$$

b) Hence, or otherwise, calculate the betweenness of the i th node from the end of a "line graph" of n nodes, i.e., n nodes in a row like this:

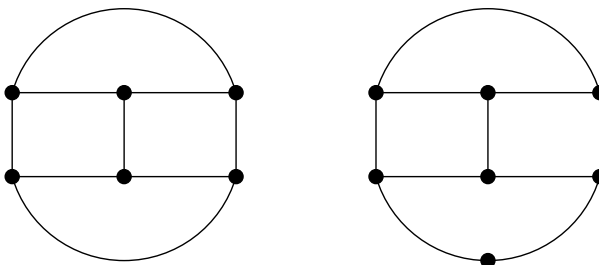


7.8 Consider these three networks:

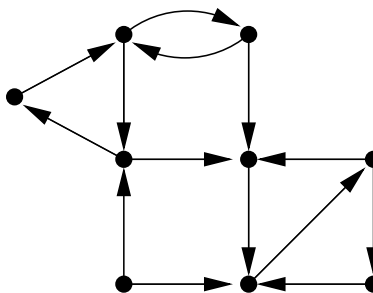


- a) Find a 3-core in the first network.
 - b) What is the reciprocity of the second network?
 - c) What is the cosine similarity of nodes A and B in the third network?
- 7.9 Consider the following networks.

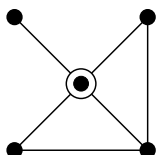
a) Find a 3-core in these two networks or state that there is none:



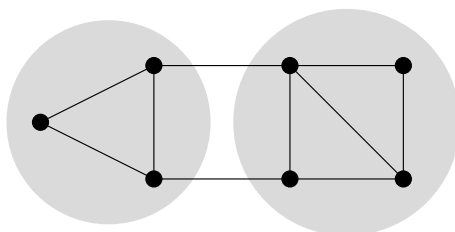
b) Find all the strongly connected components in this graph:



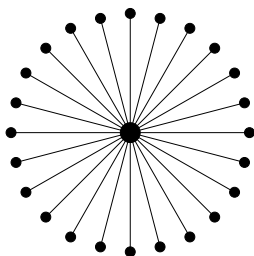
c) Calculate the local clustering coefficient of each node in this network:



d) The two circled groups of nodes in the following network represent people from Mars (on the left) and people from Venus (on the right). What is the modularity Q of the network with respect to planet of origin?



e) A “star graph” consists of a single central node and $n - 1$ other nodes connected to it thus:

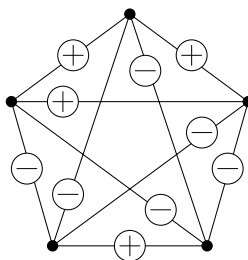


What is the (unnormalized) betweenness centrality, Eq. (7.24), of the central node as a function of n ?

7.10 What is the difference between a 3-component and a 3-core? Draw a small network that has one 3-core but two 3-components.

7.11 Among all pairs of nodes in a directed network that are connected by an edge, half are connected in only one direction and the rest are connected in both directions. What is the reciprocity of the network?

7.12 In this network + and - indicate pairs of people who like each other or don't, respectively:



- a) Is the network structurally balanced and why?
- b) Is it clusterable and, if so, what are the clusters?

7.13 Construct a proof of Harary's balance theorem (Section 7.5) as follows. Suppose that, given a connected, undirected, signed network—a network with a single component in which every edge is either a positive or negative—we color nodes with two different colors, starting from any node we please and working outward, such that nodes at opposite ends of a positive edge are the same color and nodes at opposite ends of a negative edge are different colors. This coloring process will fail if we go around a loop, returning to a previously colored node, and the rules would require us to give that node the opposite color to the one it already has. Show that this happens if, and only if, the loop in question has an odd number of negative edges around it. Hence, argue that in a network whose loops all have an even number of negative edges, the coloring process is always possible and therefore the network is clusterable in the sense of Harary's theorem.

7.14 Construct a proof of Davis's theorem that an undirected, signed network is clusterable if and only if it contains no loops with exactly one negative edge (see Section 7.5). To do this, consider the network formed by the positive edges only and imagine coloring each component of this network a different color, then adding the negative edges back into the network one by one. Argue that the network is not clusterable if and only if any negative edge falls between two nodes of the same color, and hence prove Davis's theorem.

7.15 Demonstrate that for node similarities σ_{ij} defined according to Eq. (7.50) the sum $\sum_j \sigma_{ij}$ gives the PageRank of node i divided by the degree k_i .

7.16 Consider the following two studies:

- a) In a survey of heterosexual couples in the city of San Francisco, Catania *et al.* [97] recorded, among other things, the ethnicity of their interviewees and calculated the fraction of couples whose members were from each possible pairing of ethnic groups. The fractions were as follows:

		Women				Total
		Black	Hispanic	White	Other	
Men	Black	0.258	0.016	0.035	0.013	0.323
	Hispanic	0.012	0.157	0.058	0.019	0.247
	White	0.013	0.023	0.306	0.035	0.377
	Other	0.005	0.007	0.024	0.016	0.053
Total		0.289	0.204	0.423	0.084	

Assuming the couples interviewed to be a representative sample of the edges in the undirected network of relationships for the community studied, and treating the nodes as being of four types—black, Hispanic, white, and other—calculate the numbers e_r and a_r that appear in Eq. (7.58) for each type. Hence, calculate the modularity of the network with respect to ethnicity.

- b) A 2016 article on the website FiveThirtyEight.com described the results of a study by Eitan Hersh and Yair Ghitza of political alignment among couples in the United States. Hersh and Ghitza estimated the fractions of opposite-sex partners with each combination of major-party alignment (which in the US means Democratic, Independent, or Republican) to be as follows:

		Women			Total
		Democrat	Independent	Republican	
Men	Democrat	0.25	0.04	0.03	0.32
	Independent	0.06	0.15	0.05	0.26
	Republican	0.06	0.05	0.30	0.41
	Total	0.37	0.24	0.38	

Assuming these results to be representative of the network of relationships, calculate the modularity of the network with respect to political persuasion.

- c) From your results, what do you conclude about homophily in these two studies?

CHAPTER 8

COMPUTER ALGORITHMS

An introduction to methods for performing network calculations on a computer, including data structures for storing networks and algorithms for a number of standard network problems

IN THE preceding chapters of this book we have introduced various types of networks encountered in scientific study and the basic theoretical tools used to describe and quantify them. In practice, most applications of the methods we have described, and indeed most analysis involved in the contemporary study of networks, is performed using computers. In the early days of network analysis in the first part of the twentieth century, calculations were performed by hand, partly out of necessity, since computers were rare, slow, and difficult to use, but also because the networks studied were typically quite small, consisting of perhaps just a few dozen nodes or even less. These days, by contrast, networks with thousands or even millions of nodes are not uncommon. Gathering and analyzing the data for networks like these is only possible because of the advent of fast and widely available computing.

Some network calculations are simple enough that it is obvious how one would get a computer to carry them out, but many are not and performing them efficiently requires careful consideration and thoughtful programming. Even just storing a network on a computer requires some thought, since there are several ways of doing it and the choice of which to use can make a substantial difference to the performance of subsequent calculations.

In this chapter we discuss some of the techniques used for performing

network calculations on computers. A good understanding of the material introduced here will form a solid foundation if and when you want to perform your own calculations with network data.

To begin with in this chapter we describe some simple but important ideas about the running time of algorithms and about data structures for the storage of networks. This material forms a foundation for the remainder of the chapter, in which we describe a selection of fundamental network algorithms, including many of the classics of the field, such as algorithms for calculating centrality indices, finding components, and calculating shortest paths.

Our discussion of computer algorithms does not end with this chapter, however. Building on the concepts we develop here, we will in later chapters of the book introduce a range of other, more specialized algorithms, as the need arises, including for instance algorithms for generating randomized networks in Chapter 12, algorithms for community detection in Chapter 14, and algorithms for network percolation in Chapter 15.

Understanding the content of this chapter does not require that you know how to program a computer. We will not, for instance, discuss particular programming languages. However, some experience with programming will certainly help in understanding the material, and the reader who has none will in any case probably not have much use for the methods we describe. Conversely, readers who already have a thorough knowledge of computer algorithms may find some of the material here too basic for them, particularly the material on run times and data structures (Sections 8.2 and 8.3). Such readers should feel free to skip material as appropriate and move quickly onto the possibly less familiar subject matter of Sections 8.4 onward. For very advanced readers who are already familiar with all the material covered here and who wish to go into the subject in greater detail, we recommend the books by Cormen *et al.* [122], which is a general computer science text on algorithms, and by Ahuja *et al.* [9], which is specifically on network algorithms.

8.1 SOFTWARE FOR NETWORK ANALYSIS AND VISUALIZATION

Before we leap into the study of algorithms, a few words of advice are worthwhile. Many of the standard algorithms for the study of networks are already available for use in the form of professional network analysis software packages. These packages are often of very high quality, produced by skilled and knowledgeable programmers, and if they are adequate for your needs then there is no reason not to use them. Writing and debugging your own software can take hours or days, and there is little reason to expend that time when someone else has already done it for you. Table 8.1 lists some of the most

Name	Availability	Platform	Description
Gephi	Free	WML	Interactive network analysis and visualization
Pajek	Free	W	Interactive social network analysis and visualization
InFlow	Commercial	W	Interactive social network analysis and visualization
UCINET	Commercial	W	Interactive social network analysis
yEd	Free	WML	Interactive visualization
Visone	Free	WL	Interactive visualization
Graphviz	Free	WML	Visualization
NetworkX	Free	WML	Python library for network analysis and visualization
JUNG	Free	WML	Java library for network analysis and visualization
igraph	Free	WML	C/R/Python libraries for network analysis

Table 8.1: A selection of software implementing common network algorithms. Platforms are Microsoft Windows (W), Apple Macintosh (M), and Linux (L). Most Linux programs also run under Unix and Unix-like systems such as BSD, and many Windows programs can run on Mac and Linux systems using emulation software.

widely used current software packages for the analysis of network data along with a brief description of what they do. The present author, for instance, has made considerable use of Graphviz, Pajek, Gephi, NetworkX, and yEd, all of which provide useful features that could save you a lot of time in your work. Some other network calculations can be performed using standard mathematical software such as R, Matlab, or Mathematica, and again there is no reason not to make use of these resources if they are adequate for the particular task before you.

That said there are still some excellent reasons for studying network algorithms and computer methods. First of all, even when you are making use of pre-packaged software to do your calculations, it helps greatly if you understand how the algorithms work and what the software is doing. Much time can be wasted when people fail to understand how a program works or misunderstand the kinds of answers the program can give them. Furthermore, if you are going to undertake a substantial amount of work using network data, you will sooner or later find that you need to do something that cannot be done with standard software and you will have to write some programs of your own.

Relying on pre-packaged software can create other problems too. In particular, it has a tendency to steer researchers towards investigating questions that can be answered using the software they have and away from other potentially interesting questions that would require writing new software. In this way, software packages can in effect shape the research agenda of the field, which is the reverse of the way things should work. Good research decides the interest-

ing questions first and then goes in search of answers. Research that restricts itself only to the questions it already knows how to answer will be narrowly focused indeed.

By following the developments in this and the succeeding chapters, and, if you wish, reading further in the algorithms literature, you give yourself the opportunity to pursue whatever network questions are of interest to you, without having to rely on others to produce software to tackle those questions.

8.2 RUNNING TIME AND COMPUTATIONAL COMPLEXITY

Before we look at exactly how network algorithms work, there is an important issue we need to tackle, that of *computational complexity*. If you have programmed computers before, you may well have had the experience of writing a program to perform a particular calculation and setting it running, only to find that it is still running an hour or even a day later. Performing a quick back-of-the-envelope calculation, you discover to your dismay that your program is going to take a thousand years to finish, and hence that it is basically useless.

The concept of computational complexity (or just “complexity” for short) is essentially a more formal version of back-of-the-envelope calculations like this, and is useful precisely because it helps us to avoid wasting our energy on programs that will not finish running in any reasonable amount of time. By considering the complexity of an algorithm before we even start to write a computer program, we can be sure we are writing one that will actually do the job.

Computational complexity is a measure of the running time of a computer algorithm, as a function of the size of the problem it is tackling. Consider a simple example: how long does it take to find the largest number in a list of n numbers? Assuming the numbers are not given to us in some special order (such as largest first), then there is no quicker way to find the largest than simply to go through the whole list, number by number, keeping a running record of the largest one we have seen, until we get to the end.

This is a very simple example of a computer algorithm. We could use it, for instance, to find the node in a network that has the highest degree. The algorithm consists of a set of steps. At each step we examine the next number in the list and ask whether it is larger than the largest we have seen so far. If it is, it becomes the new largest-number-so-far. Otherwise, nothing happens and we move on to the next step.

Now here is the crucial point: in the *worst possible case* the most work we will have to do for this algorithm is on each step to (1) compare the next number in the list with our previous record holder and (2) replace the previous record

holder with the new number. That is, the largest amount of work we have to do happens when every number is bigger than all those before it.

In this case the amount of work we do is the same on every step and hence the total time taken to complete the algorithm, its running time, is just $n\tau$, where τ is the time taken on each individual step. If we are lucky, the actual time taken may be less than this, but it will never be more. Thus, we say that the running time or *time complexity* of this algorithm is of order n , or just $O(n)$ for short.

Technically, the notation $O(n)$ means that the running time varies as a constant times n or less, to leading order in n .¹ We say “to leading-order” because it is possible that there may be contributions to the running time that increase with n more slowly than this leading-order term. For instance, there might be some initial start-up time for the algorithm, such as time taken initializing variables, that is a constant independent of n . We would denote this time as being $O(1)$, i.e., a constant times 1. By convention, however, one drops such sub-leading terms when citing the complexity of an algorithm, because if n is large enough that the running time of the program becomes a serious issue, then the sub-leading terms will usually be small enough by comparison with the leading ones that they can be safely neglected.² Thus, by the conventional definitions, the time complexity of our simple largest-number algorithm is just $O(n)$.

The computational complexity of an algorithm is an indication of how the algorithm’s running time scales with the size of its input. In our example, the input to the algorithm is the list of numbers and the size of that input is the length n of the list. If this algorithm were used to find the highest degree node in a network, for instance, the size of the input would be the number of nodes. In many of the network algorithms we will look at this will be the case—the number of nodes n will be the important parameter we consider. In other cases, the important parameter will be the number of edges m in the network, while in others still we will need both m and n to fully specify the size of the input—there could be different parts to an algorithm, for instance, that operate separately on the nodes and the edges, so that the total running time depends on both. Thus, for example, we will see in Section 8.5 that the algorithm known as breadth-first search, which is used for finding shortest paths between nodes, has a computational complexity of $O(m) + O(n)$ for a network with m edges

¹If we wish to say that the running time is exactly proportional to n , we can use the notation $\Theta(n)$.

²There are occasional instances where this is not true, so it is worth just bearing in mind the possibility of sub-leading terms.

and n nodes, meaning that it runs in time $am + bn$ where a and b are constants, or quicker. Very often one writes this, in shorthand and slightly sloppy form, as $O(m + n)$. This latter notation is not meant to imply that the constants in front of m and n are the same.

In a lot of networks research we are concerned with sparse networks, and particularly with the networks we called “extremely sparse” in the discussion of Section 6.10.1, meaning those for which the average degree $c = 2m/n$ tends to a constant as n becomes large. This means that m increases in proportion to n , which in turn implies that $O(m + n)$ is equivalent to $O(n)$. In this common case, therefore, we can drop the m from our notation.

The importance of computational complexity lies in its use for estimating the actual running time of real algorithms. Suppose, for example, that we wish to run the breadth-first search algorithm mentioned earlier on a network with a million nodes and ten million edges. Knowing that the algorithm has time complexity $O(m + n)$, we could start out with a small test-run of the program on a network with $n = 1000$ nodes, say, and $m = 10\,000$ edges. Often we artificially create small networks just for the purposes of such tests.

Perhaps we find that the program finishes in one second on the test network. We can then scale up this result knowing that the running time varies as $am + bn$. On the full network with $n = 1\,000\,000$ and $m = 10\,000\,000$ both n and m are a thousand times larger than on the test network, so the program should take about a thousand times longer to finish, i.e., a thousand seconds or about a quarter of an hour. Armed with this information we can safely start our program working on the larger problem and step out for a cup of coffee or a phone call while we wait for it to finish.

Conversely, suppose we had an algorithm with computational complexity $O(n^4)$. That means that if we increase the number of nodes n in our network by a factor of a thousand the running time will increase by a trillion. If, for instance, we try the algorithm on a small test network and find that it takes a second again, then a network a thousand times larger would take a trillion seconds, which is around 30 000 years. In this case, we would certainly abandon the calculation, or at least look for a faster algorithm that can complete it in reasonable time.

Finding the computational complexity of an algorithm, generating test networks, performing short runs, and doing scaling calculations of this type all require some work—additional work on top of the work of developing and programming the computer algorithm in the first place. Nonetheless, this extra work is well worth the effort involved and one should always perform this type of analysis, at least in some rough manner, before embarking on any major numerical calculation. Computational complexity will be one of our principal

concerns in the discussions of algorithms in this and succeeding chapters. In practice, an algorithm is useless for all but the smallest of networks if its running time scales poorly with the size of a network. As a general rule, any algorithm that scales with system size as $O(n^3)$ or greater is too slow for use on large networks, although such algorithms might still find some use for the smaller cases. In the world of computer science, where many researchers have devoted their entire careers to the invention of new algorithms for solving particular problems, the calculation of the computational complexity of an algorithm is a primary goal—often *the* primary goal—of research. Plenty of papers are published whose sole contribution is to provide a calculation of the complexity of some algorithm.

It is worth mentioning that calculations of the running time of algorithms based on their complexity, as above, do not always give completely accurate answers. We have mentioned already that standard measures of time complexity neglect sub-leading contributions to the run time, which may introduce inaccuracies in practical situations. But in addition there are, for technical reasons, cases where the behavior of the running time is poorer than a simple scaling argument would suggest. For instance, in calculations on networks it is important that the entire network fit in the main memory (RAM) of our computer if the algorithm is to run quickly. If the network is so large that at least part of it must be stored on a disk or some other slow form of storage, then the performance of the algorithm may be substantially hindered.³ Even if the entire network fits in the main memory, there may be additional space required for the operation of the algorithm, and that must fit in the memory too. Also, not all kinds of memory are equally fast. Modern computers have a small amount of extra-fast “cache” memory that the computer can use for storing small amounts of frequently used data. If all or most of the data for a calculation fit in the cache, then the program will run significantly faster than if they do not.

There are also cases where a program will perform better than its time complexity might indicate. In particular, the complexity is usually calculated by considering the behavior of the program in the worst case. But for some programs the worst-case behavior is relatively rare, occurring only for particularly unlucky values of the program inputs, and the typical behavior is significantly better than the worst case. For such programs the complexity can give an

³There are whole subfields in computer science devoted to the development of algorithms that run quickly even when part of the data is stored on a slow disk. Usually such algorithms work by reordering operations so that many operations can be performed on the same data, stored in the main memory, before swapping those data for others on the disk.

unreasonably pessimistic estimate of running time.

Despite these caveats, however, computational complexity is still a useful guide to overall program performance and an indispensable tool in the computer analysis of large networks.

8.3 STORING NETWORK DATA

The first task of most programs that work with network data is to read the data, usually from a computer file, and store it in some form in the memory of the computer. Network data can be stored in files in any of a large number of different formats, some standard, some not, but typically a file includes an entry with information about each node or about each edge, or sometimes both. However, it is the way the data are stored in the computer memory after they are read from the file that has the biggest effect on the running of a program. As we will see, different choices about how to store the data can make a substantial difference to both the speed of a program and the amount of memory it uses. Here we discuss some of the commonest ways to store network data in computer memory.

The first step in representing a network in a computer is to label the nodes so that each can be uniquely identified. The most common way of doing this is to give each a numeric label, usually an integer, just as we have been doing in our mathematical treatment of networks in previous chapters. It usually does not matter which node gets assigned which number—the purpose of the numbers is only to provide unique labels for identifying the nodes. In the simplest case we number the n nodes of a network by the consecutive integers $i = 1 \dots n$, although in some cases we might use non-consecutive integers or start the numbering from a different point. For instance, in some programming languages, including C, Python, and Java, it is conventional for numbering to start at zero and go up to $n - 1$. Most, though not all, file formats for storing networks already specify integer labels for nodes, in which case we often just use those labels. For those that don't, one typically labels nodes consecutively in the order they are read from the file. In what follows, we will assume that nodes are numbered $1 \dots n$.

Often the nodes in a network have other notations or values attached to them in addition to their integer labels. The nodes in a social network, for instance, might have names; nodes on the World Wide Web might have URLs; nodes on the Internet might have IP addresses or AS numbers. Nodes could also have properties like age, capacity, or weight represented by additional numbers, integer or not. All of these other notations and values can be stored straightforwardly in the memory of the computer by defining an array of a suit-

able type with n elements, one for each node, and filling it with the appropriate values in order. For example, we might have an array of n text strings to store the names of the individuals in a social network and another array of integers to store their ages in years.

Having devised a suitable scheme for labeling the nodes and storing their properties, we need a way to represent the edges of the network. This is where things get more complicated.

8.3.1 THE ADJACENCY MATRIX

In most of the mathematical developments of previous chapters we have represented networks by their adjacency matrix \mathbf{A} . The adjacency matrix also provides one of the simplest ways to represent a network on a computer. Most computer languages provide two-dimensional arrays that can be used to store a matrix directly in memory. An array of integers can be used to store an adjacency matrix if the matrix consists only of integers, as it does for unweighted networks and multigraphs. An array of floating-point numbers would be needed for an adjacency matrix that has non-integer elements, as occurs in some weighted networks.

Storing a network in the form of an adjacency matrix is convenient in many ways. Most of the formulas and calculations described in this book are written in terms of adjacency matrices. So if we have that matrix stored in our computer it is usually a trivial matter to turn the formulas into computer code and calculate the corresponding quantities.

The adjacency matrix makes other operations straightforward too. For instance, if one wishes to add or remove an edge between a given pair of nodes, this can be achieved quickly and easily with an adjacency matrix. To add an edge between nodes i and j in an unweighted network one simply increases the ij th element of the adjacency matrix by one. To remove an edge between the same nodes one decreases the element by one. These operations take a constant amount of time regardless of the size of the network, so their computational complexity is $O(1)$. Similarly if we want to test whether there is an edge between a given pair of nodes i and j we need only inspect the value of the appropriate matrix element, which can also be done in $O(1)$ time.

Undirected networks give a slight twist to the situation since they are represented by symmetric matrices. If we want to add an undirected edge between nodes i and j , then in principle we should increase both the ij th and ji th elements of the adjacency matrix by one, but in practice this is a waste of time. A better approach is to update only elements in the upper triangle of the matrix and leave the lower one empty, knowing that its correct value is just the mirror

See Section 6.2 for an introduction to the adjacency matrix.

Adjacency matrices for weighted networks are discussed in Section 6.3.

image of the upper triangle. (For directed networks, which are represented by asymmetric adjacency matrices, this issue does not arise—the full matrix, both the upper and lower triangles, is used to store the structure of the network.) To put this another way, in an undirected network we should only update elements A_{ij} of the adjacency matrix for which $i < j$. For instance, if we wish to create an edge between node 2 and node 1, this means in principle that we want to increase both the (2, 1) element and the (1, 2) element of the adjacency matrix by one. But, since we are only updating elements with $i < j$, we would increase only the (1, 2) element and leave the other alone.

Taking this idea one step further, we could decline to store the lower triangle of the adjacency matrix in memory at all. If we are not going to update it, why waste memory storing it? Unfortunately, dropping the lower triangle of the matrix makes our remaining matrix triangular itself, and most computer languages don't provide triangular arrays. One can, with a certain amount of work, arrange to store triangular sets of quantities using, for instance, the dynamic memory allocation facilities provided by languages like C and Java, but this is only worth the extra effort if memory space is the limiting factor in performing your calculation.

The adjacency matrix is not always a convenient representation, however. It is cumbersome if, for instance, we want to run quickly through the neighbors of a particular node, at least on a sparse network. The neighbors of node i are denoted by non-zero elements in the i th row of the adjacency matrix and to find them all we would have to go through all the elements of the row one by one looking for those that are non-zero. This takes time $O(n)$ (since that is the length of the row), which could be a lot of time in a large network, and yet on a sparse network most of that time is wasted, because most of the elements in the adjacency matrix are zero. As we will see in this chapter, many network algorithms do indeed require us to find all neighbors of a node, often repeatedly, and for such algorithms the adjacency matrix is not an ideal tool.

The computational complexity of the network operations discussed here for an adjacency matrix is summarized in Table 8.2.

Another disadvantage of the adjacency matrix representation is that for the common case of a sparse network it makes inefficient use of computer memory. In a network in which most elements of the adjacency matrix are zero, most of the memory occupied by the matrix is used to store those zeros. As we will shortly see, there is an alternative representation known as an adjacency list that avoids storing the zeros and thereby takes up much less space.⁴

We could equally well store the edges in the lower triangle of the matrix and neglect the upper triangle. Either choice works fine.

For networks in which self-edges are allowed, we would use the diagonal elements as well, so we would update elements with $i \leq j$ —see Section 6.2.

⁴One advantage of the adjacency matrix is that the amount of space it consumes is independent of the number of edges in the network (though it still depends on the number of nodes). As we will

Operation	Adjacency matrix	Adjacency list
Insert	$O(1)$	$O(1)$
Delete	$O(1)$	$O(m/n)$
Find	$O(1)$	$O(m/n)$
Enumerate	$O(n)$	$O(m/n)$

Table 8.2: The time complexity of four basic network operations. The leading-order time complexity of four operations when carried out with adjacency matrix and adjacency list representations of a network of n nodes and m edges. The operations are adding an edge to the network (insert), removing an edge from the network (delete), testing whether a given pair of nodes are connected by an edge (find), and listing the neighbors of a given node (enumerate).

It is a simple matter to work out how much memory is consumed in storing the adjacency matrix of a network. The matrix has n^2 elements. If each of them is an integer (which requires 4 bytes for its storage on most modern computers) then the entire matrix will take $4n^2$ bytes. At the time of writing, a typical computer has about 10^{10} bytes of RAM (10 GB), and hence the largest network that can be stored in adjacency matrix format satisfies $4n^2 = 10^{10}$, or $n = 50\,000$. This is not nearly large enough to store the largest networks we have encountered, such as large subsets of the Web or large social networks, and is not even big enough for some of the medium-sized ones.

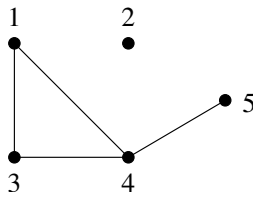
The disadvantages of the adjacency matrix representation described here apply primarily to sparse networks. If one is interested in dense networks—those in which a significant fraction of all possible edges are present—then the adjacency matrix format may be appropriate. It will still use a lot of memory in such cases, but so will any data format, since there is simply a lot of information that needs to be stored, so the advantages of other formats are less significant. The adjacency matrix may also be a good choice if you are only interested in relatively small networks. For instance, the social network analysis package UCINET, which is targeted primarily at sociologists working with smaller networks, uses the adjacency matrix format exclusively. A lot of current research on networks, however, is focused on larger data sets, and for these another representation is needed.

see in the next section, adjacency lists use varying amounts of memory, even for networks with the same number of nodes, depending on how many edges there are. In calculations where edges are frequently added or removed it may be convenient—and increase the speed of our algorithms—to have the size of our data structures remain constant, although this advantage must be weighed against the substantial space savings of using the adjacency list.

8.3.2 THE ADJACENCY LIST

The most common alternative to storing the complete adjacency matrix of a network is to use an *adjacency list*. The adjacency list is in fact probably the most widely used method for storing networks on a computer.

An adjacency list is actually not just a single list but a set of lists, one for each node i . Each list contains the labels of the other nodes to which i is connected. Consider, for example, this small undirected network:



which would be represented by this adjacency list:

Node	Neighbors
1	3, 4
2	
3	4, 1
4	5, 1, 3
5	4

An adjacency list can be stored in a series of integer arrays, one for each node, or as a two-dimensional array with one row for each node.⁵ It is common to also store somewhere the degree of each node, so that we know how many entries there are in the list of neighbors for each node. This can be done using a separate array of n integers. Note also that there is usually no requirement that the neighbors of a node in the adjacency list appear in numerical order. Normally they are allowed to appear in any order.

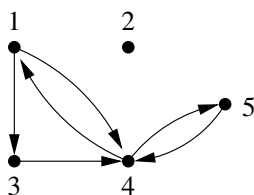
Looking at the example adjacency list above, you will notice that each edge appears twice. For instance, node 3 is listed as a neighbor of node 1 and node 1 is also listed as a neighbor of node 3. The adjacency list for a network with m edges therefore consists of $2m$ integers. This is much better than the n^2 integers

⁵Note that the number of entries in the list of neighbors for a node is equal to the degree of the node and can vary from one node to another. It may even be zero. Most modern computer languages, including C and Java, allow the creation of two-dimensional matrices with rows having varying numbers of elements, making it straightforward to store adjacency lists in a memory-efficient way. Some older languages, like FORTRAN 77, do not allow this, making things more difficult.

used to store the full adjacency matrix.⁶ For instance, on a computer where each integer occupies 4 bytes of memory, a network with $n = 10\,000$ nodes and $m = 100\,000$ edges would occupy 800 kB in adjacency list form, as opposed to 400 MB in matrix format. The double storage of the edges is somewhat wasteful—we could save an additional factor of two if we only stored each edge once. However, the double storage turns out to have some advantages, making our algorithms substantially faster and easier to program in many cases, and these benefits are normally worth the extra cost in terms of space. In these days of cheap memory, not many networks are large enough that the space required to store an adjacency list presents a serious problem.

An adjacency list can store networks with multiedges or self-edges. A multi-edge is represented by multiple identical entries in the list of neighbors of a node, all pointing to the same adjacent node. A self-edge is represented by an entry identifying a node as its own neighbor. In fact, a self-edge is most correctly represented by *two* such entries in the list, so that the total number of entries in the list is still equal to the degree of the node. (Recall that a self-edge adds 2 to the degree of the node it is connected to.)

The example adjacency list above is for an undirected network, but adjacency lists can be used with directed networks as well. Consider, for instance, this network:



which can be represented by the adjacency list⁷

⁶Note that the amount of memory used is now a function of m rather than n . For algorithms in which edges are added or removed from a network during the course of a calculation this means that the size of the adjacency list can change, which can complicate the programming and potentially slow down the calculation. Normally, however, this added complication is not enough to outweigh the considerable benefits of the adjacency list format.

⁷Indeed, the adjacency list for an undirected network could be viewed as a special case of the directed adjacency list for a network in which each undirected edge is replaced by two directed ones, one in each direction. It takes only a moment to convince oneself that this results precisely in the sort of double representation of each edge that we saw in the undirected case.

Node	Outgoing edges
1	3, 4
2	
3	4
4	5, 1
5	4

Here we have listed only the outgoing edges for each node. Since each edge is outgoing from some node, this approach is guaranteed to capture every edge in the network, but each edge now appears only once in the adjacency list, not twice as in the undirected case.

Alternatively, we could represent the same network by listing the ingoing edges for each node thus:

Node	Incoming edges
1	4
2	
3	1
4	3, 1, 5
5	4

In principle these two representations contain the same information. Both include all the edges and either of them can be constructed from a knowledge of the other. When creating computer programs, however, the crucial point is to have the information you need for your calculations easily available, so that the program runs quickly. Different calculations require different information and some might need ingoing edges while others need outgoing ones. The choice of which adjacency list to use thus depends on the particular calculations being performed. Some calculations might even require both ingoing and outgoing edges, in which case we could create a double adjacency list like this:

Node	Incoming edges	Outgoing edges
1	4	3, 4
2		
3	1	4
4	3, 1, 5	5, 1
5	4	4

Note that, as in the undirected case considered above, this double adjacency list stores each edge twice, once as an incoming edge and once as an outgoing one, and is thus in some respects wasteful of space, although not to an extent that is often a problem.

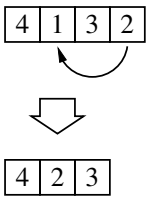
As with the adjacency matrix, it is important to ask how fast our calculations will run if we store our network as an adjacency list. Will they run at a reasonable pace? If the answer is no, then the adjacency list is not a useful representation, no matter what its other advantages may be.

Consider the undirected case⁸ and the four basic network operations that we considered previously for the adjacency matrix, addition and removal of edges, finding whether an edge is present, and enumeration of all edges connected to a node—see Table 8.2.

We can add an edge to our adjacency list very quickly: to add an edge between nodes i and j we need only add one new entry each to the ends of the neighbor lists for nodes i and j , which takes time $O(1)$.

Finding or removing an edge is a little harder. To find whether an edge exists between nodes i and j we need to go through the list of neighbors of i to see whether j appears in that list, or vice versa. Since the list of neighbors is normally in no particular order, there is no quicker way of doing this than simply going through the entire list step by step from the beginning. In the worst case we will have to check all elements to find our edge or confirm that it does not exist, and on average⁹ this will take time of order the mean number c of elements in the list, which is given by the mean degree $c = 2m/n$ (Eq. 6.15). Thus the “find” operation takes time $O(m/n)$ for a network in adjacency list form. This is a bit slower than the same operation using an adjacency matrix, which takes time $O(1)$ —see Section 8.3.1. On a sparse network with constant mean degree (see Sections 6.10 and 8.2), $O(m/n)$ is equivalent to $O(1)$, so technically the complexity of the adjacency list is as good as that of the adjacency matrix, but in practice the former will be slower than the latter by a constant factor which could become large if the average degree is large.

Removing an edge involves first finding it, which takes time $O(m/n)$, and then deleting it. The deletion operation can be achieved in $O(1)$ time by simply moving the last entry in the list of neighbors to overwrite the entry for the deleted edge and decreasing the degree of the node by one (see figure). (If the edge we are deleting is the last element, then we need do nothing other than decreasing the degree by one.) Thus the leading-order running time for the



The element “1” is deleted from a list by moving the last element “2” to overwrite it.

⁸The answers are essentially the same in the directed case. The demonstration is left as an exercise.

⁹We are thus calculating a sort of “average worst-case” behavior, allowing for the worst case in which we must look through the entire list, but then averaging that worst case over many different lists. This is a reasonable approach because almost all of the algorithms we will be considering do many successive “find” operations during a single run, but it does mean that we are technically not computing the complexity of the absolute worst case situation.

edge removal operation is $O(m/n)$.

However, the adjacency list really comes into its own when we need to run quickly through the neighbors of a node, a common operation in many network calculations, as discussed in Section 8.3.1. We can do this very easily by simply running through the stored list of neighbors for the node in question, which takes time proportional to the number of neighbors, which on average is $c = 2m/n$. The leading-order time complexity of the operation is thus $O(m/n)$, much better than the $O(n)$ of the adjacency matrix for the same operation.

The computational complexity of operations on the adjacency list is summarized in Table 8.2.

8.3.3 OTHER NETWORK REPRESENTATIONS

We have discussed two common ways of representing network data in the memory of a computer: the adjacency matrix and adjacency list. These are the representations you are most likely to use if you write your own programs, but there are others that are also worth knowing about.

Hybrid matrix/list representations: The adjacency matrix and adjacency list both have advantages and disadvantages and neither is optimal in all cases. In the best of all possible worlds, we would like a data structure that can insert, delete, and find edges in $O(1)$ time and enumerate the $O(m/n)$ neighbors (on average) of a given node in $O(m/n)$ time, but neither the adjacency matrix nor the adjacency list can do this. It is possible to create a representation that *can* do this, however, if we are willing to sacrifice memory space. We describe two ways of doing this.

One approach is to make a hybrid representation that consists of an adjacency matrix *and* an adjacency list. Non-zero elements in the adjacency matrix, those corresponding to edges, are accompanied by pointers that point to the corresponding elements in the adjacency list. Then we can find whether an edge exists between a specified pair of nodes in $O(1)$ time using the adjacency matrix as usual. And we can enumerate the neighbors of a node in $O(m/n)$ time using the adjacency list. We can add an edge in $O(1)$ time since both matrix and list allow this anyway (Table 8.2). And we can delete an edge in $O(1)$ time by first locating it in the adjacency matrix and setting the corresponding element to zero, then following the pointers to the relevant elements of the adjacency list and deleting those too by moving the last element of the list to fill their place.

In terms of time complexity, i.e., scaling of run time with network size, this

hybrid data structure is optimal.¹⁰ Its main disadvantage is that it uses even more memory than the ordinary adjacency matrix, and hence is suitable only for relatively small networks, up to a few tens of thousands of nodes on a typical computer at the time of writing. If this is not an issue in your case, however, and speed is, then this hybrid representation may be worth considering.

Adjacency lists stored in data structures other than arrays: Perhaps a more satisfactory approach than the hybrid matrix/list is to represent the network using an adjacency list, but to record the neighbors of each node in a different kind of data structure. Instead of using an array, we can use any one of a number of alternative structures that allow for faster finding and removal of elements, usually at the expense of greater programming complexity or memory use. An example of such a data structure is a balanced tree, such as a *AVL tree* or *red-black tree* [122]. These are standard data structures, which we will not describe in detail here, that allow one to add, find, and remove elements in time proportional to the log of the number of elements in the tree, and enumerate elements in time proportional to the number of elements. Thus, if we use a separate tree to store the list of neighbors of each node, we can perform the addition, removal, and finding of edges¹¹ in time $O(\log(m/n))$ and their enumeration in time $O(m/n)$. While this is not truly optimal (an optimal approach would use $O(1)$ time for the addition, removal, and find operations) it is still pretty good. The logarithm is a slowly increasing function of its argument, so $\log(m/n)$ is typically a small number, and for the common case of a sparse network with $m \propto n$ it is actually constant.

In fact, there are other data structures that do even better, allowing one to add, remove, and find elements in constant time and enumerate all elements in time proportional to the total number of elements. In particular, the data structure called a *hash table* can achieve this performance on average. (That is, the time for the add, remove, and find operations is constant on average but individual operations may take longer or shorter times, depending on particular details of the data stored.) An adjacency list in which the set of

¹⁰It does place some overhead on the edge addition and deletion operations, meaning the complexity is still $O(1)$ but the operations take a constant factor longer to complete, since we have to update both adjacency matrix and list, where normally we would only have to update one or the other. Whether this makes a significant difference to the running time of a program will depend on the particular algorithm under consideration.

¹¹The time to add, remove, or find an edge connected to node i goes as $\log k_i$ and hence the average time goes as $\langle \log k \rangle$, where $\langle \dots \rangle$ denotes the average over all nodes. This is not necessarily the same as $\log \langle k \rangle = \log(2m/n)$. However, the result as quoted is still correct because $\langle \log k \rangle \leq \log \langle k \rangle$ always, so the running time is of order $\log(m/n)$ or less, which is precisely the meaning of the notation $O(\log(m/n))$.

neighbors of each node is stored in a separate hash table can therefore perform the addition, removal, find, and enumerate operations in time $O(1)$, $O(1)$, $O(1)$, and $O(m/n)$, respectively, which are optimal. The primary disadvantage of hash tables is that they use more memory space than the simple array-based adjacency list. How much memory they use depends on how fast we want them to work. A hash table has a parameter called the *load factor* that controls the payoff between speed and memory use, whose value one can vary depending on how much one cares about these two issues. The total amount of memory used is always proportional to the number of values stored in the table, but the constant of proportionality varies with the load factor. For a typical load factor a hash table might, for example, use twice as much memory to store the same values as a simple array. Thus the total amount of memory needed to store a network using hash tables is proportional to the number of edges m in the network, which is much better than the hybrid adjacency matrix/list approach of the previous section, for which the total amount of memory is $O(n^2 + m)$, which is usually much larger. (The n^2 is for the adjacency matrix and the m is for the adjacency list.)

For these reasons, an adjacency list stored in a set of hash tables may well be the best possible network representation, at least for some applications. It gives optimal performance for the four basic network operations we have considered, and does so with a level of memory use that still allows one to store very large networks—potentially up to billions of nodes and edges on a typical computer at the time of writing. If one can live with the (relatively modest) extra memory use and the more complicated programming required to use hash tables rather than arrays, then this is an approach well worth considering.

A detailed description of how hash tables work can be found, for example, in the book by Cormen *et al.* [122]. However, many modern computer languages have hash tables already built in, so you don't have to program them yourself. They are a standard feature in the Java language, for example, and also in Python (where they are called "dicts"). Other languages, such as C and C++, provide hash tables through the use of external libraries.

Representations with variables on edges: In some networks the edges have values, weights, or labels on them. One can store additional properties like these using simple variants of the adjacency matrix or adjacency list representations. For instance, if edges come in several types we could use a range of integer values in the elements of the adjacency matrix to indicate edge types. If there are several different variables associated with each edge, as there are for instance in some social network studies, then we could use several different matrices, one for each variable, or a single matrix whose elements are themselves arrays

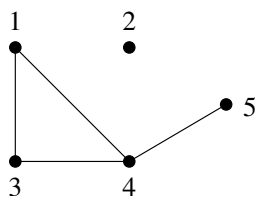
of values or other more complicated objects. Similarly, with an adjacency list we could replace the elements of the list with objects that contain all the details of the edges they correspond to, an approach that works whether those objects are stored in an array, a tree, a hash table, or some other data structure.

However, these representations can be wasteful or clumsy. The matrix method can waste huge amounts of memory on all the empty matrix elements for edges that don't exist. And the adjacency list (for an undirected network) contains two entries for each edge, both of which would have to be updated every time we modify the properties of an edge.

In some cases, therefore, it is worthwhile creating additional data structures to store the properties of the edges separately. For instance, one might use an array of m elements, one for each edge. This array can be linked to the main representation of the network structure: with an adjacency list, for instance, we could store in each entry in the list a pointer to the corresponding element in the array of edge data. Then we can immediately find the properties of any edge we encounter in the main adjacency list. Similarly, each entry in the array of edge data could include pointers to the elements in the adjacency list that correspond to the edge in question. This would allow us to go through the array of edge data looking for edges with some particular property and, for example, delete them.

Some standard libraries for storing and manipulating networks, such as the Python NetworkX library, provide automatic mechanisms for separately storing node and edge properties in this way. For complex network programming projects you may wish to consider using such a library. There will be an extra investment of time required to learn how to use it, but in the long run it may well save you time by doing a lot of the trickier steps of the programming for you. There is no reason to avoid using standard libraries if they help you get things done faster.

Edge lists: One very simple representation of a network that we have not yet mentioned is the *edge list*. This is simply a list of the pairs of nodes that are connected by edges. Consider, for instance, this network, which we saw previously in Section 8.3.2:



The edge list representation would be $(1, 3), (4, 1), (4, 3), (4, 5)$. The order of the edges is usually not important in an edge list, nor is the order of the nodes in the node pairs.

The edge list is a convenient and space-efficient way to store the structure of a network, and furthermore allows us easily to associate properties with edges—we can simply store those properties along with the corresponding pairs of labels. It is not, however, a very good representation if we wish to store properties of nodes. Indeed, an edge list doesn't explicitly list nodes at all, so there is no way to tell that a node even exists if it is not connected to any edges. Node 2 in the network above is an example: it doesn't appear in the edge list because it has no edges. On the other hand, this problem and the problem of storing node properties can be remedied easily enough by creating a separate list of nodes and the data associated with them.

Even with this modification, however, the edge list is a poor format for storing network data in computer memory in most cases. It does not, for instance, allow us to determine quickly whether a particular edge exists—we would have to go through the whole list of edges to answer that question. And, crucially, it does not allow us easily to enumerate the neighbors of a given node, an operation that is central to many algorithms. For these reasons, the edge list is hardly ever used as a format for the representation of a network in the memory of a computer.

Where it does find use is in file formats for networks. Being a fairly compact representation, edge lists are often used as a way to store network structure in computer files on a disk drive or other storage medium. When we wish to perform calculations on these networks we must read the file and convert its contents into a more suitable form for calculation, such as an adjacency list. This, however, is simple: we create an empty network in the memory of our computer, one with no edges, in the format of our choice—an empty adjacency list, for instance. Then we run through the edges stored in the edge list and add them one by one to the network. Since the operation of adding an edge can be accomplished quickly in all of the formats we have considered, this process normally does not take a significant amount of time. When it is finished, we have a complete copy of the network stored in the memory of the computer in our desired format, and we are ready to start our computations.

8.4 ALGORITHMS FOR BASIC NETWORK QUANTITIES

Armed with the computational tools introduced in previous sections for storing and manipulating network data, we turn now to a discussion of specific algorithms for network calculations. We start with a brief discussion of some

simple algorithms for calculating quantities such as degrees and clustering coefficients, then spend the remainder of the chapter looking in detail at the more complex algorithms used for calculating shortest paths, betweenness, maximum flows, and other non-local quantities.

8.4.1 DEGREES

Many network quantities are easy to calculate and require only the simplest of algorithms, algorithms that are little more than translations into computer code of the definitions of the quantities in question. Nonetheless, it is worth looking at these algorithms at least briefly, for two reasons. First, there is in many cases more than one method for calculating a quantity, and some methods may be faster than others. It pays to evaluate one's algorithm at least momentarily before writing a computer program, to make sure one is going about the calculation in the most sensible manner. Second, it is worthwhile to calculate the computational complexity of even the simplest algorithm, so that one can make an estimate of how long a computation will take to finish—see Section 8.2. Even simple algorithms can take a long time to run.

As an example, consider the degree of a node, one of the most fundamental and important of network quantities. Normally, degrees are very simple to calculate. In fact, if a network is stored in the form of an adjacency list then, as described in Section 8.3.2, we normally maintain an array containing the degree of each node so that we know how many entries there are in the list of its neighbors. That means that finding the degree of any particular node is a simple matter of looking it up in this array, which takes $O(1)$ time.

If the network is stored as an adjacency matrix, then the calculation takes longer. Calculating the degree of a node i in this case involves going through all elements of the i th row of the adjacency matrix and counting the number that are non-zero. Since there are n elements in each row of the matrix, where n is the number of nodes in the network, the calculation takes time $O(n)$, far longer than for the adjacency list. If one needed to find the degrees of nodes frequently during the course of a larger calculation using an adjacency matrix, it might make good sense to calculate the degree of each node once and for all and store the results for later easy retrieval in a separate array.

Other quantities related to degree are similarly straightforward to calculate. Take, for example, the correlation coefficient r for node degrees, Eq. (10.27), which measures assortative mixing by degree. The correlation coefficient can be calculated using Eq. (10.28) and the sums defined in Eqs. (10.29) and (10.30). Given the degrees of all nodes, the sum in Eq. (10.29) takes time $O(m)$ to evaluate, where m is the number of edges in the network, and the sums in (10.30)

each take time $O(n)$, so the total time required to calculate r is $O(m + n)$. As mentioned in Section 8.2, we are often concerned with what we called “extremely sparse” networks, those in which the mean degree remains constant as the network gets larger, i.e., networks in which $m \propto n$. In such networks $O(m + n)$ is equivalent to $O(n)$ and the time taken to calculate r just scales as the number of nodes. On the other hand, if the network is dense, meaning that $m \propto n^2$, then $O(m)$ is equivalent to $O(n^2)$, which is considerably worse.

8.4.2 CLUSTERING COEFFICIENTS

The calculation of clustering coefficients is only slightly more complicated than the calculation of degrees. To see how it works, we start by calculating the local clustering coefficient, Eq. (7.29), for a single node i on an undirected network:

$$C_i = \frac{\text{(number of pairs of neighbors of } i \text{ that are connected)}}{\text{(number of pairs of neighbors of } i)}}. \quad (8.1)$$

Calculating the numerator involves going through every pair of distinct neighbors of node i and counting how many are connected. We need only consider each pair once, which we can do conveniently by restricting ourselves to pairs u, v for which $u < v$. For each pair we determine whether an edge exists between them, which is done in different ways depending on the representation used for the network (see Section 8.3), and count up the number of such edges. Then we divide the result by the number of pairs, which is just $\frac{1}{2}k_i(k_i - 1)$, where k_i is the degree of the node as usual.

To calculate the overall clustering coefficient for the entire network, which is given by

$$C = \frac{\text{(number of triangles)} \times 3}{\text{(number of connected triples)}}, \quad (8.2)$$

(see Eq. (7.28)), we extend the same calculation to the whole network. That is, for every node we consider each pair of neighbors u, v with $u < v$ and find whether they are connected by an edge.¹² We add up the total number of such edges over all nodes and then divide by the number of connected triples, which is $\sum_i \frac{1}{2}k_i(k_i - 1)$.

This algorithm is simple and straightforward, a direct implementation of Eq. (8.2), but some interesting issues nonetheless come up when we consider its running time. Most of the effort in the algorithm is taken up with counting

¹²Note that this calculation automatically accounts for the factor of three appearing in the numerator of Eq. (8.2), since each triangle is counted three times, once from the point of view of each of the three nodes it connects.

the connections between pairs of neighboring nodes. We must check for the presence of an edge between each such pair in the entire network and hence the total number of checks we must perform is

$$\sum_i \frac{1}{2} k_i (k_i - 1) = \frac{1}{2} n (\langle k^2 \rangle - \langle k \rangle), \quad (8.3)$$

where

$$\langle k \rangle = \frac{1}{n} \sum_i k_i, \quad \langle k^2 \rangle = \frac{1}{n} \sum_i k_i^2, \quad (8.4)$$

are the mean and mean-square degree for the network. (We previously denoted the mean degree by c , but we use the alternate notation $\langle k \rangle$ here for clarity, and to highlight the distinction between the mean and the mean square.)

The interesting point here is that Eq. (8.3) depends in a non-trivial way on the degrees in our network. The running times of other algorithms we have seen so far have depended on the number of nodes n and the number of edges m , and sometimes on the ratio m/n , which is proportional to the mean degree $\langle k \rangle = 2m/n$. For the clustering coefficient, however, we see that the amount of work we must do, and hence also the running time, depends not only on n and the mean degree but also on the mean square.

In Section 10.4 we will look at “scale-free networks,” a special type of network whose node degrees follow a power-law distribution. For such networks, as we will see, the mean degree is typically well behaved but the mean square formally diverges, which implies that it will take an *infinite* amount of time to evaluate the clustering coefficient. Even if we use a representation of the network, such as an adjacency matrix, that allows us to check for the presence of an edge rapidly in time $O(1)$, the number of checks we must do still diverges and the time taken will be infinite.

To be fair, as discussed in Section 10.4.2, the second moment does not actually become infinite for any finite network. But it may become very large, which means that the calculation of the clustering coefficient will be slow.

These difficulties are specific to the case of scale-free networks. In other cases there is usually no problem calculating the clustering coefficient quickly. Some alternative algorithms have been proposed for calculating approximate values of the clustering coefficient rapidly, such as the algorithms of Schank and Wagner [423] and Tsourakakis [448], and these may be worth considering if you need to perform calculations on very large networks.

8.5 SHORTEST PATHS AND BREADTH-FIRST SEARCH

We now turn to some more complex algorithms, such as algorithms for calculating distances, flows, and cut sets between nodes. Our discussion of each of these algorithms will have three parts. Two are the description of the algorithm and the analysis of its running time. But now we also include a proof that each algorithm actually performs the calculation it claims to. In the case of algorithms like those for degrees and clustering coefficients such proofs are unnecessary; the algorithms are simple translations of the equations for the corresponding quantities. As we move on to more complex algorithms, however, it will become much less obvious why they give the results they do, and some sort of proof, either formal or informal, is often necessary to convince ourselves that we are doing the calculation correctly.

The first algorithm we look at is the standard algorithm for finding distances between nodes in a network, which is called *breadth-first search*. One run of the breadth-first search algorithm finds the shortest distance from a single starting node s to every other node in the same component of the network. With only minor modifications it can also find the path one must take to realize each shortest distance, and if there is more than one shortest path it can find all of them. It works on both directed and undirected networks, although our description will focus on the undirected case.

In some situations we want to know only the shortest distance between a single pair of nodes s, t , in which case we could just use breadth-first search to calculate distances to all nodes, then throw away all of the results except the one we want. This, however, is rather wasteful and there is a variant of breadth-first search that can do the calculation more quickly—see Section 8.5.4.

In physics, breadth-first search is sometimes called the “burning algorithm.”

8.5.1 DESCRIPTION OF THE BREADTH-FIRST SEARCH ALGORITHM

Breadth-first search finds the shortest network distance from a given starting node s to every other node in the same component as s . The basic principle is illustrated in Fig. 8.1. The algorithm proceeds through a series of rounds. In the first round we find all the neighbors of s , which by definition have distance 1 from s . In the next round we find all the neighbors of the neighbors, which have distance 2, then the neighbors of the neighbors of the neighbors, which have distance 3, and so on. At every round we go one step further out from node s , until there are no nodes left to explore.

A more formal definition of breadth-first search is as follows. The algorithm works by assigning numbers to nodes one by one, indicating their distance from s . Initially node s is assigned distance 0 and all other nodes are

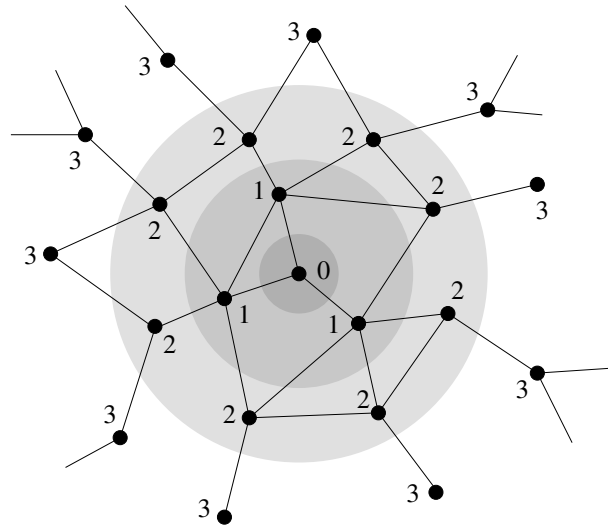


Figure 8.1: Breadth-first search. A breadth-first search starts at a given node, which by definition has distance 0 from itself, and grows outward in layers or rounds, labeling nodes with their distance from the starting node as it goes. The nodes explored in the first round, which are the immediate neighbors of the starting node, have distance 1. The neighbors of those neighbors have distance 2, and so forth.

unassigned. Then, for each integer value of d from 0 upward, we locate all nodes with distance d , find their neighbors, and mark any that currently have no assigned distance as having distance $d + 1$. The algorithm ends when there are no nodes with the next value of d .

It may be obvious to you that this algorithm does what it claims to and calculates all shortest distances from s . But just to be thorough we should prove it. The proof is by induction. We make the claim that after the d th round of the algorithm all nodes up to and including distance d from node s —and only those nodes—have been assigned their correct distances. Assuming this to be true, we note the following two facts: (1) For a node at distance d , every neighbor that has not been assigned a distance by the end of the d th round must have distance $d + 1$, and (2) every node at distance $d + 1$ is a neighbor of at least one node with distance d . These two statements may be obvious to you, but if not, Exercise 8.6 on page 273 gives you an opportunity to prove them for yourself.

On the next round of the algorithm we locate all nodes at distance d and mark all of their currently unassigned neighbors as having distance $d + 1$. This process finds all nodes with distance $d + 1$ (because of (2) above) and only those nodes (because of (1) above). Hence after one more round of the algorithm we have correctly assigned distances to all nodes up to distance $d + 1$ from s , and only those nodes. Thus, if we repeat the process for increasing values of d , all reachable nodes will eventually be marked with their correct distances, once the value of d reaches the largest distance of any node from s .

It remains only to provide a “base case” to get the induction started. The base case here is the case for $d = 0$. There is only one node with distance 0, the node s itself, which is correctly assigned its distance when the algorithm starts, all other nodes being unassigned. With this base case the proof of the algorithm is complete.

As a by-product of its operation, breadth-first search also finds the component to which node s belongs, since the algorithm assigns distances only to those nodes that can be reached from previously assigned ones, meaning only nodes in the component. At the end of the calculation the set of nodes with assigned distances thus corresponds to the component containing s and all other nodes are unassigned. Indeed, breadth-first search is the algorithm of choice for finding components in networks.

8.5.2 A NAIVE IMPLEMENTATION

Let us consider how we would implement breadth-first search on a computer. The simplest approach (but not, as we will see, the best) would go something like this. We create an array of n integer elements to store the distance of each node from the source node s , and initially set the distance of node s from itself to zero while all other nodes have unknown distance from s . Unknown distances could be indicated, for instance, by setting the corresponding element of the array to -1 , or some other value that could never occur in reality.

We also create a distance variable d to keep track of where we are in the breadth-first search process and set its value initially to zero. Then, following the prescription of the previous section, we do the following:

1. Find all nodes that are distance d from s by going through the distance array, element by element. If there are no nodes with distance d the algorithm ends.
2. Find all the neighbors of those nodes and check each one to see whether its distance from s is unknown (denoted, for example, by -1 in the distance array). If its distance is unknown, assign it distance $d + 1$.
3. Increase the value of d by 1.
4. Repeat from step 1.

When the algorithm is finished we are left with an array that contains the distance from s to every node in the same component as s (and every node in every other component has no assigned distance).

How long does this algorithm take to run? Let us go through the operations performed by the algorithm in turn. First of all we have to set up the array of distances, giving each of the n elements its appropriate initial value. We

spend a constant amount of time assigning a value to each element, so overall we spend $O(n)$ time setting up the array.

For the algorithm proper, on each round we go through all n nodes looking for those with distance d , which takes time $O(n)$. If the total number of rounds is r then the overall time spend on this part of the algorithm will be $O(rn)$. (We will discuss the value of r in a moment.)

When we do come across a node with distance d , we must pause at that node and spend some additional effort checking each of its neighbors to see whether their distances are unknown and assigning them distance $d + 1$ if they are. If we assume that the network is stored in adjacency list format (see Section 8.3.2) then we can go through the neighbors of a node in time $O(m/n)$ on average, and during the whole course of the algorithm we pause like this at each node exactly once, so that the total extra time we spend checking neighbors of nodes is $n \times O(m/n) = O(m)$. (Strictly speaking, we only have to pause at each node in the component containing s , but in the worst case where the network is composed of a single component this means all n nodes.)

Thus, the total running time of the algorithm, including the time taken for set-up, is $O(n + rn + m)$.

And what is the value of the quantity r ? The value of r is the total number of rounds of the algorithm, which is equal to the maximum distance from our source node s to any other node. In the worst case, this distance is equal to the diameter of the network (the largest distance from any node to any other—see Section 6.11.1) and the worst-case diameter for a network with n nodes is $n - 1$, which is realized when the network is just a chain of nodes strung one after another in a line. Thus, in the worst case our algorithm will have running time $O(m + n^2)$ (where we are keeping only the leading-order terms in the expression).

This is pessimistic, however. As we will see in Sections 10.2 and 11.7, the diameter of most networks increases only as $\log n$, in which case our algorithm runs in time $O(m + n \log n)$ to leading order, which is significantly better. This may be a moot point, however, since we can do better still if we use a little cunning in the implementation of our algorithm.

8.5.3 A BETTER IMPLEMENTATION

The most time-consuming part of the algorithm described in Section 8.5.2 is the part where we go through the nodes of the network to find those that are exactly distance d from the starting node s . Since this operation involves checking all n nodes, only a small fraction of which may be at distance d , it often wastes a lot of time.

the element pointed to by the read pointer and increase the pointer by one.

3. Find the distance d for that node by looking in the distance array.
4. Go through the neighbors of the node in turn and look up their distances in the distance array as well. If a neighbor has a known distance, leave it alone. If it has an unknown distance, mark it as having distance $d + 1$, store its node label in the queue array, in the element pointed to by the write pointer, and increase the write pointer by one.
5. Repeat from step 2.

Note the test applied in step 2: if the read pointer points to the same element as the write pointer, then there is no node to be read from the queue (since the write pointer always points to an empty element). Thus this test tells us when there are no further nodes waiting to have their neighbors investigated.

The algorithm reads all the nodes with distance d from the queue array one after another and uses them to find all the nodes with distance $d + 1$. Thus, all nodes with the same distance appear one after another in the queue array, with the nodes of distance $d + 1$ immediately after those of distance d . Furthermore, each node appears in the queue array at most once. A node may be a neighbor of more than one other, but it is assigned a distance and put in the queue only on the first occasion on which it is encountered. If it is encountered again, its distance is then known rather than unknown, and hence it is not again added to the queue. Of course, a node may not appear in the queue array at all if it is never reached by the breadth-first search process, i.e., if it belongs to a different component from s .

Thus the queue does exactly what we wanted it to: it stores all nodes with distance $d + 1$ for us so that we have the list handy on the next round of the algorithm. This spares us from having to search through the network for these nodes and so saves us a lot of time. In all other respects the algorithm works exactly as before and gives the same answers.

How long does this implementation of the algorithm take to run? Again we first have to set up the array of distances, and we also now have to set up the n -element queue array. Together these operations take time $O(n)$. Then, for each element in the queue, which means for each of the nodes in the same component as s , we do the following operations: we run through its neighbors, of which there are $O(m/n)$ on average, and either calculate their distance and add them to the queue, or do nothing if their distance is already known. Either way the operations take $O(1)$ time. Thus, for each node in the component, of which there are in the worst case n , we spend time $O(m/n)$ and hence we require overall at most a time $n \times O(m/n) = O(m)$ to complete the algorithm for all n nodes.

Thus, including the time to set up the arrays, the whole algorithm takes time $O(m + n)$, which is better than the $O(m + n \log n)$ of the naive implementation (Section 8.5.2). For the common case of a sparse network with $m \propto n$, $O(m + n)$ is equivalent to $O(n)$ and our algorithm runs in time proportional to the number of nodes. This is optimal, since the algorithm is calculating the distance to all n nodes from the source node s , and one cannot assign n numbers to n array elements in less than n time.

On a sparse network, therefore, the breadth-first search algorithm does as well as we can hope for in finding the distances from a single node to all others, and indeed it is the fastest known algorithm for performing this operation.

On the other hand, for a dense network where $m \propto n^2$, we have a running time of $O(n^2)$.

8.5.4 VARIANTS OF BREADTH-FIRST SEARCH

There are a number of variants of breadth-first search that merit a mention. First, one might wish to calculate the shortest distance between only a single pair of nodes s and t , rather than between s and all others. As mentioned at the start of Section 8.5, one way to do this is just to use breadth-first search to calculate the distance from s to every node and then throw away all the results except for the one we want. If we wanted to be a little cleverer about it, we could save some time by running the breadth-first search only until we find and assign a distance to node t , then stop. There is no need to continue the calculation beyond this point.

This still wastes a lot of effort, however, and there is a simple variant of the algorithm that can do much better. The trick is to perform *two* breadth-first searches at the same time, starting from the two nodes of interest, s and t , as shown in Fig. 8.2. In successive rounds of the algorithm we find and label all nodes at distance 1 from either s or t , then all nodes at distance 2, and so forth. The algorithm ends the first time one of the two breadth-first searches assigns a distance to a node that has already been assigned a distance by the other search. The total distance from s to t is then the sum of these two distances (i.e., it is equal to the distance from s to the node plus the distance from the node to t).

This *two-source breadth-first search* is typically significantly faster than the standard single-source version because it has to probe fewer nodes before it gets the answer it needs. In the single-source version, the algorithm starts from node s and runs until it finds node t . If we are lucky, node t might be the first node we look at and the algorithm will end after just one step. In the worst case, on the other hand, we find t last and have to look through all n nodes before we get to it. On average, therefore, we expect that we will have to look at about $\frac{1}{2}n$ nodes before the algorithm stops. For each of these nodes we must

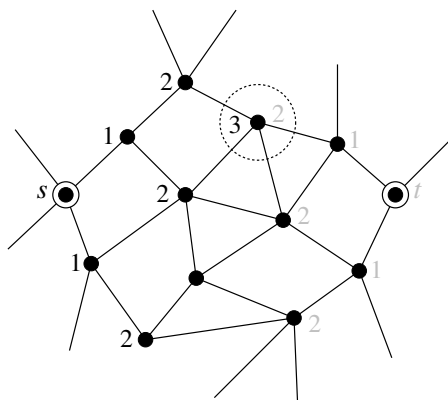


Figure 8.2: The two-source breadth-first search algorithm. To find the shortest distance between s and t we start two simultaneous breadth-first searches to calculate the distances from s to other nodes (numbers in black) and from t to other nodes (numbers in gray). On the first round of both searches we label nodes at distance 1 from their respective starting points. On the second round we label nodes at distance 2. On the third round we start labeling nodes at distance 3, but when we come to label the circled node at the top as having distance 3 we discover that it has already been previously labeled as having distance 2 from t . At this point the algorithm stops and the shortest distance from s to t is the sum $3 + 2 = 5$.

look at its neighbors, of which there are $O(m/n)$ on average, so the total work done is $O(m)$ again, as with standard breadth-first search (plus any time taken setting up the arrays).

The number of nodes examined in the two-source version of the algorithm is harder to pin down because it depends on the details of the network. However, as we will see in Section 11.7, it is typically the case that the number of nodes reached by a breadth-first search grows roughly exponentially with the distance d probed, i.e., as c^d for some constant c . Thus, a single-source search between an s and t a distance d apart will examine about c^d nodes. As we have said, however, the number of nodes examined by a single-source search is about $\frac{1}{2}n$ on average, so $c^d \simeq \frac{1}{2}n$ for a typical s, t . The two-source algorithm, on the other hand, stops when each of its two searches has probed out to distance about $\frac{1}{2}d$, which means they each examine about $c^{d/2} \simeq \sqrt{n/2}$ nodes. Then the total number of nodes examined by both searches is twice this figure, or $\sqrt{2n}$. Multiplying by the time $O(m/n)$ to examine the neighbors of each node, the total running time is then $O(m/\sqrt{n})$, a factor of \sqrt{n} better than the single-source algorithm. For a network with a million nodes, for example, we'd expect the

two-source algorithm to be about a thousand times faster on average than the single-source version.¹³

Conversely, we sometimes want to calculate the shortest distance between every pair of nodes in an entire network. This we can do by performing a standard breadth-first search starting at each node in the network in turn. The total running time for this “all-pairs shortest distance” calculation is $n \times O(m + n) = O(n(m + n))$, or $O(n^2)$ on a sparse network with $m \propto n$. As with the standard breadth-first search, this is optimal in the sense that we are calculating $O(n^2)$ quantities in $O(n^2)$ time, which is the best we can hope for. The same calculation can also be used to find the diameter of a network (Section 6.11.1), which is the longest distance between any pair of nodes. In general, there is no faster way to calculate the diameter than to perform a breadth-first search starting from every node and record the largest distance observed in any of them.

The closeness centrality of Section 7.1.6 can also be calculated in a straightforward manner using breadth-first search. Recall that the closeness centrality of a node is defined as the inverse of the mean distance from that node to all others in the same component. Since our breadth-first search already calculates distances to all others in a component, we need only go through the distance array and calculate the average of all assigned distances, then take the inverse, to get the closeness value for a node. Again the running time is $O(m + n)$. The variant closeness defined in terms of the harmonic mean in Eq. (7.22) can also be calculated, in the same running time, by a similar method.

8.5.5 FINDING SHORTEST PATHS

The breadth-first search algorithm as we have described it finds the shortest distance from a node s to all others in the same component of the network. It

¹³This doesn’t count the $O(n)$ time needed to initialize the arrays, which could be a problem on a sparse network with $m \propto n$: while the running time of the main algorithm on a sparse network is $O(m/\sqrt{n}) = O(\sqrt{n})$, the total running time including initialization would still be $O(n)$ to leading order, which is as bad as the single-source algorithm. One can get around this issue by storing the distances not in an array but in a different kind of data structure, such as a hash table. Instead of storing distances for all nodes, one stores only the assigned distances and not the unassigned ones. The hash table allows us to do this but still to quickly find distances or determine that a node does not have an assigned distance. Because we store only the assigned distances, the hash table is much quicker to set up at the start of the algorithm—it is initially empty and can be set up in time $O(1)$ —and the overall running time of the two-source algorithm becomes $O(\sqrt{n})$ again. A hash table is more complicated than a simple array and using it does slow down the algorithm a bit, which is why hash tables are not normally used for the standard one-source algorithm. But for the two-source case it can be invaluable. See Cormen *et al.* [122] for a discussion of hash tables.

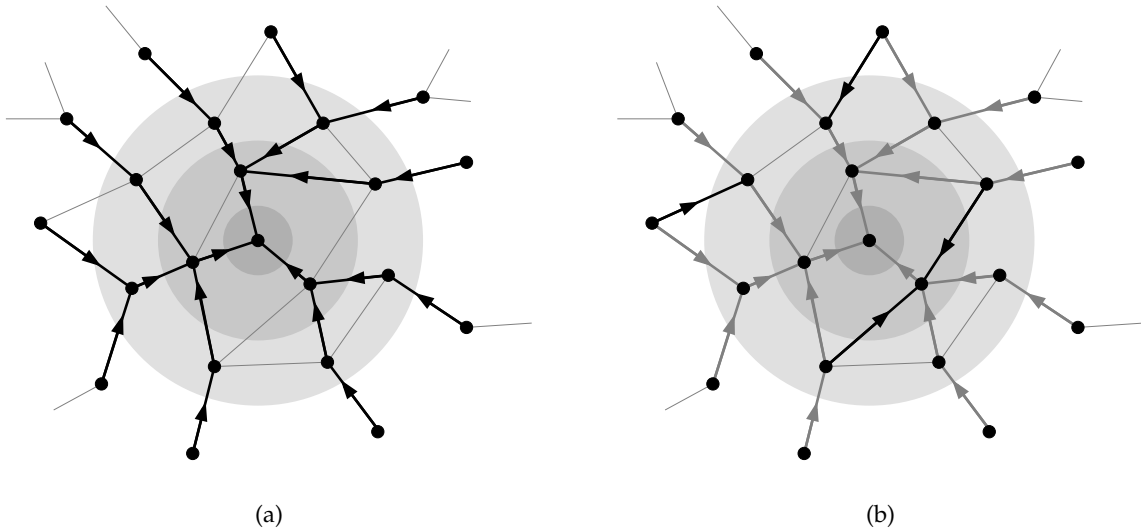


Figure 8.3: Shortest-path trees. (a) A simple shortest-path tree for the network of Fig. 8.1. Each node has an arrow or pointer leading to its “predecessor,” the node from which it was reached during the breadth-first search process. By following pointers from any node we can find a shortest path to the starting node in the center. (b) The full shortest-path tree (which is actually not a tree now but a directed acyclic graph) contains extra pointers that allow us to reconstruct all possible shortest paths.

does not tell us the particular path or paths by which that shortest distance is realized. With only a relatively small modification of the algorithm, however, we can calculate the paths as well. The idea is as follows.

We conduct our breadth-first search algorithm as before, searching outward from the starting node s . The algorithm repeatedly pulls a node from the queue and examines its neighbors as described in Section 8.5.3. But now every time the neighbor j of some node i turns out to be a previously unseen node, one with no assigned distance, we not only assign j a distance, we also create a pointer from node j to its *predecessor*, node i . The pointer can be simply an integer associated with node j that contains i 's node label. For instance, we could create a new integer array of n elements, one for each node, to store the pointers.

We continue this process until the breadth-first search finishes, and when we are done we have an array of distances as usual, but also we have a pointer from each node to its predecessor. Figure 8.3a shows a depiction of the pointers for the same network as in Fig. 8.1—they are represented as bold arrows from

nodes to their predecessors. Armed with these we can, starting at any node, now follow a pointer to the node's predecessor, then to *that* node's predecessor, and so on until we eventually get all the way back to s . The end result is a complete path back to s , and moreover this is a shortest path, since every pointer from a node at distance d points, by definition, to a node at distance $d - 1$, so that there must be exactly d steps before we reach s .

Note that the pointers, the arrows in Fig. 8.3a, form a tree—there is exactly one path from every node to s and hence no loops. This tree is sometimes called the *shortest-path tree* for node s : it is a compact representation of shortest paths from every node to s .

The extra step of creating the pointers can be accomplished quickly—in $O(1)$ time if we store the pointers in a simple array as described above. Thus the overall running time of the algorithm is the same as for standard breadth-first search: it takes time $O(m + n)$ to find all distances from s and construct the shortest-path tree.

The algorithm does have one shortcoming, which is that it only finds *one* shortest path to each node. As pointed out in Section 6.11.1, a pair of nodes may have more than one shortest path between them (see Fig. 6.12). Another slight modification of the algorithm allows us to deal with this case.

Multiple shortest paths exist between any node and the starting node s if the path to s splits in two or more directions somewhere along its length. This occurs if there is a node j at some point along the path, say at distance $d + 1$ from s , that has more than one neighbor at distance d , i.e., more than one predecessor—see Fig. 8.3b. We can record this circumstance in our shortest-path tree by adding more than one pointer from j —one to each of the predecessor nodes. To do this we modify our algorithm as follows.

We perform the breadth-first search starting from s as before, and add a pointer from each newly found node to the node from which it was reached. But we also add an extra step. If, in the process of examining the neighbors of a node i that has distance d from the source node, we discover a neighbor j that already has an assigned distance, and that distance is $d + 1$, then we know that a path of length $d + 1$ has already been found to j , but we also know that *another* path of length $d + 1$ must exist via the current node i . So we add an extra pointer from j to i . This makes the shortest-path tree no longer actually a tree, although, being a little sloppy with their nomenclature, people often call it a “shortest-path tree” anyway. In any case, the algorithm now gives exactly what we want. When it is finished running, the shortest-path “tree” contains all shortest paths from every node in the component to the source node s . See

Fig. 8.3b again.¹⁴

We can employ essentially the same method in the case where we want to find the shortest path between only a single pair of nodes s, t , but using the two-source version of breadth-first search introduced in Section 8.5.4. As described in that section, we perform two simultaneous breadth-first searches, searching outward from nodes s and t , stopping when the same node gets assigned a distance by both searches. Then we construct the shortest-path trees for both searches and the shortest path between s and t is given by the path from s to the common node plus the path from the common node to t . As described in Section 8.5.4, the calculation has a typical running time of $O(m/\sqrt{n})$, which is a factor of \sqrt{n} faster than using a single-source search.

If there is more than one shortest path between s and t and we want to find all of them, then things are a little more complicated. In that case we again carry out our two breadth-first searches until they meet at a common node, but now, to ensure we find all shortest paths, we must continue until we have completed the current round of each search, finding all further nodes that have the same distances from s and t as the common node. By doing this we ensure that if there are other common nodes we will find those also. Then we again construct the two shortest-path trees, and the shortest paths between s and t are those that pass through any of the common nodes. It is possible that there is more than one path that goes through a single common node, if there is more than one path from s to the common node, or more than one from the common node to t , or both—we must include every combination of the first and second halves of the path to find all shortest paths.

8.5.6 BETWEENNESS CENTRALITY

In Section 7.1.7 we described betweenness centrality, a widely used centrality index that measures the extent to which a node in a network lies on the paths between other nodes. The betweenness centrality of node v is the number of shortest paths between pairs of nodes s, t that pass through v . (Sometimes it is normalized to be the fraction of such paths, rather than the total number. The difference is only a multiplicative constant—see Section 7.1.7.) Given that we have a method for finding the shortest path (or paths) between any two nodes (Section 8.5.5), we can with only a little more work now create an algorithm for

¹⁴Storing the pointers is a little more complicated now. A simple integer array with one element per node will not suffice. One could use a two-dimensional array, or better still a dynamic data structure such as a linked list or hash table, whose size adjusts automatically to accommodate the pointers as they are added.

calculating betweenness.

The simplest way to calculate betweenness would be to implement the definition of the measure directly: find the shortest path between s and t , as described in Section 8.5.5 (assuming such a path exists), and then work our way along that path checking the nodes it passes through to see whether the node v we are interested in lies among them. Repeating this process for every distinct pair s, t , we can then count the total number of paths that pass through v . (Things are slightly more complicated for the case in which a pair of nodes are connected by more than one shortest path, but let us ignore this complication for the moment—we will come to it soon.)

This is certainly a correct algorithm and it would work, but it is also inefficient. As we have seen, the best way to find the shortest path between a single pair of nodes s, t is to use the two-source version of breadth-first search (Section 8.5.4), which takes time $O(m/\sqrt{n})$. Since there are $\frac{1}{2}n(n-1)$ distinct pairs of nodes s, t , finding the shortest paths between all of them would take $O(n^{3/2}m)$ time, or $O(n^{5/2})$ in the common case of a sparse network for which $m \propto n$. This is prohibitively slow: it might work for networks up to a few thousand nodes, but for the larger networks that are typical of modern network studies the calculation would not be feasible.

We can, however, do a lot better if we make use of some results from previous sections. First, recall that the standard (single-source) breadth-first search algorithm can find paths between a source s and all other nodes (in the same component) in time $O(m+n)$, which means, as noted in Section 8.5.4, that we can find paths between all pairs in the network in time $O(n(m+n))$, or $O(n^2)$ on a sparse network where $n \propto m$.

An improved algorithm for calculating the betweenness of a node v might thus work as follows. For each node s we perform a breadth-first search starting at that node and then construct a shortest-path tree as described in Section 8.5.5. We use that tree to trace the paths from each node back to s and make a note of how many of those paths go through v . We repeat this calculation for all s and so end up with a count of the total number of shortest paths¹⁵ that pass through v . Indeed, we can trivially extend this algorithm to calculate betweenness for *all* nodes at the same time—we simply maintain a count of the

¹⁵Note that this approach actually counts each path twice (since the path between i and j is counted once when i is considered the source node and once when j is), except for the path from each node to itself, which is counted only once (when that node is the source). This, however, is correct: the betweenness centrality, as defined in Eq. (7.24), indeed counts each path twice, except for the path from a node to itself. As mentioned in Section 7.1.7, some researchers prefer to use a definition that counts each path only once. The latter can be easily calculated from the former, however, by adding one and dividing by two.

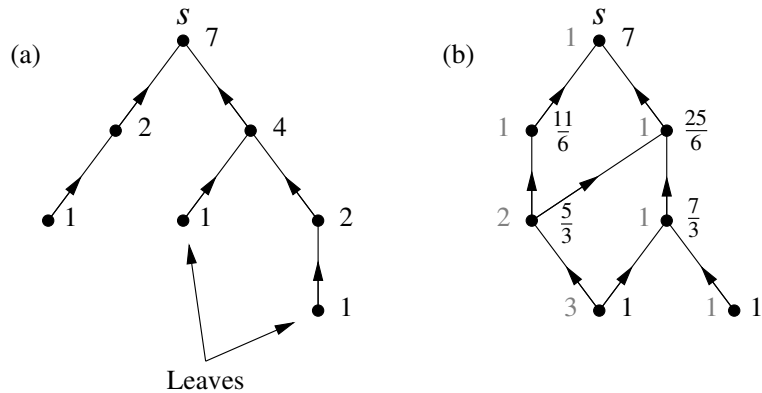


Figure 8.4: Calculation of betweenness centrality. (a) When there is only a single shortest path between a source node s (top) and all other reachable nodes, those paths necessarily form a tree, which makes the calculation of the contribution to betweenness from this set of paths particularly simple, as described in the text. (b) For cases in which there is more than one shortest path to some nodes, the calculation is more complex. First we must calculate the number of paths between the source s and each other node (numbers to the left of nodes), and then use these to weight the path counts appropriately and derive the betweenness scores (numbers to the right of nodes).

number of paths that go through every node, for example in an array.

For any given s , this algorithm will take time $O(m + n)$ to find the shortest paths. Shortest paths by definition have length less than or equal to the diameter of the network, which is typically of order $\log n$ (see Sections 10.2 and 11.7), and hence traversing the n paths from each node to s will take time $O(n \log n)$, for a leading-order running time of $O(m + n \log n)$ for each value of s . Repeating for all s , the whole algorithm will then take total time $O(n(m + n \log n))$ or $O(n^2 \log n)$ on a sparse network.

This is an improvement on our earlier $O(n^{5/2})$ algorithm, but we can do better still. It is in fact possible to cut the running time down to just $O(n(m + n))$ by exploiting the fact that many of the paths in the shortest-path tree share edges. To understand this development, consider Fig. 8.4a, which shows an example of a shortest-path tree for a breadth-first search starting at node s . For the moment we consider the simple case where there is only one shortest path between s and any other node, so that the shortest-path tree really is a tree, as depicted. We will consider the more general case of multiple shortest paths in a moment.

We use the tree to calculate a score for each node representing the number of

shortest paths passing through that node. We first find the “leaves” of the tree, i.e., those nodes such that no shortest paths from other nodes to s pass through them. In Fig. 8.4a the leaves are drawn at the bottom of the tree. We assign a score of 1 to each of these leaves—the only path to s that passes through these nodes is the one that starts there.¹⁶ Then, starting at the bottom of the tree we work upward, assigning to each node a score that is 1 plus the sum of the scores on the nodes immediately below it. That is, the number of paths through a node v is 1 for the path that starts at v itself plus the count of all paths that start below v and hence have to pass through it.

When we have worked all the way up the tree in this manner and reached node s , the scores at each node are equal to the betweenness counts for paths that end at node s . Repeating the process for all s and summing the scores, we arrive at the full betweenness scores for all paths.

In practice, the process of working our way up the tree can be accomplished by running through the nodes in order of decreasing distance from s . Conveniently, we already have a list of nodes in order of their distances, namely the entries in the queue array created by the breadth-first search process. Thus, the betweenness algorithm in practice involves running backwards through the list of nodes in the queue array and calculating the number of paths through each node as above until the beginning of the array is reached.

The amount of work involved in this process is proportional to the number of edges in the shortest-path tree, which in the worst case is just equal to the number m of edges in the network itself. The breadth-first search takes time $O(m + n)$ (as usual) and hence the total time to count paths for each source node s is also $O(m + n)$, which means the complete calculation of betweenness for all n nodes takes time $O(n(m + n))$, as promised.

In general, however, we cannot assume that there is only a single shortest path between every pair of nodes. As we saw in Section 8.5.5, there can be more than one, in which case the shortest-path “tree” is not actually a tree at all. Consider, for instance, Fig. 8.4b, in which some nodes have more than one shortest path to node s . Following the definition of betweenness in Section 7.1.7, such multiple shortest paths are given equal weights summing to 1, so that for a node pair connected by three shortest paths, for example, we give each path

¹⁶In this case we are considering the first and last nodes on a path to be members of that path. As mentioned in footnote 12 on page 174, the first and last nodes are sometimes excluded from the definition of betweenness, which means that the betweenness score of each node is smaller by an additive constant equal to twice the number of nodes in the component. If we wish to calculate betweenness according to this definition, the simplest approach is to use the algorithm described here and then subtract the additive constant from each node’s score at the end.

weight $\frac{1}{3}$. Note that some of the paths may follow the same route for part of their length, in which case those portions receive weight equal to the sum of the weights for the corresponding paths.

To correctly calculate the weights of the paths flowing through each node in a network, we first need to calculate the total number of shortest paths from each node to s . This is actually quite straightforward to do: the shortest paths from a node i to node s must pass through one or more neighbors of i and the total number of shortest paths from i to s is simply the sum of the numbers of shortest paths from each of those neighbors to s . We can calculate these sums as part of a modified breadth-first search process as follows.

Consider Fig. 8.4b and suppose we are starting at node s . We carry out the following steps:

1. Assign node s distance zero, place it in the queue, and set $d = 0$. Also assign node s a weight $w_s = 1$ (whose purpose will become clear shortly).
2. If the read and write pointers of the queue are equal, the breadth-first search is finished. Else read the next node from the queue. Call this node i .
3. From node i follow each attached edge to the node j at its other end and then do one of the following three things:
 - a) If j has not yet been assigned a distance, assign it distance $d + 1$ and weight $w_j = w_i$, then add it to the queue.
 - b) If j has already been assigned a distance and that distance is equal to $d + 1$, then the node's weight is increased by w_i , so that $w_j \leftarrow w_j + w_i$.
 - c) If j has already been assigned a distance less than $d + 1$, do nothing.
4. Increase d by 1.
5. Repeat from step 2.

The resulting weights for the example of Fig. 8.4b are shown to the left of each node in the figure. Each weight is the sum of those above it in the "tree." (It may be helpful to work through this example yourself by hand to see how the algorithm arrives at these values for the weights.)

Physically, the weight on a node i represents the number of distinct shortest paths between the source node s and i . Hence, if there is a pointer from j to i in the shortest-path tree, then the fraction of the paths to s that pass through (or start at) j and that also pass through i is given by w_i/w_j .

Thus, and finally, to calculate the contribution to betweenness from shortest paths starting at all nodes and ending at s , we carry out the following steps:

1. Find every "leaf" node t , i.e., a node such that no paths from s to other nodes go through t , and assign it a score of $x_t = 1$.

2. Starting at the bottom of the tree, work up towards s and assign to each node i a score $x_i = 1 + \sum_j x_j w_i/w_j$, where the sum is over the neighbors j immediately below node i .
3. Repeat from step 2 until node s is reached.

The resulting scores are shown to the right of each node in Fig. 8.4b. Now repeating this process for all n starting nodes s and summing the resulting scores gives us the total betweenness for all nodes.

This algorithm again takes time $O(n(m+n))$ in general or $O(n^2)$ on a sparse network, which is the best known running time for any betweenness algorithm at the time of writing, and moreover seems unlikely to be beaten by any future algorithm given that the calculation of the betweenness necessarily requires us to find shortest paths between all pairs of nodes, which operation also has time complexity $O(n(m+n))$. Indeed, even if we want to calculate the betweenness of only a single node it seems unlikely we can do better given that such a calculation still requires us to find all shortest paths.

8.6 SHORTEST PATHS IN NETWORKS WITH VARYING EDGE LENGTHS

In Section 6.3 we discussed weighted networks, networks in which the edges have values or strengths representing, for instance, the traffic capacities of connections on the Internet or the frequencies of contacts between acquaintances in a social network. In some cases the values on edges can be interpreted as lengths. These could be real lengths, such as distances along roads in a road network, or they could represent quantities that act like lengths, such as travel times for airline flights or transmission delays for packets traveling along Internet connections. In other cases they might just be approximately length-like measures: one might say, for instance, that a pair of acquaintances in a social network are twice as far apart as another pair if they see one another half as often.

Sometimes with networks such as these we would like to calculate the shortest path between two nodes taking the lengths of the edges into account. For instance, we might want to calculate the shortest driving route from A to B via a road network or we might want to calculate the route across the Internet that gets a data packet to its destination in the shortest time. (In fact, this is exactly what many Internet routers do when routing data packets.)

But now we note a crucial—and annoying—fact. The shortest path across a network when we take edge lengths into account may

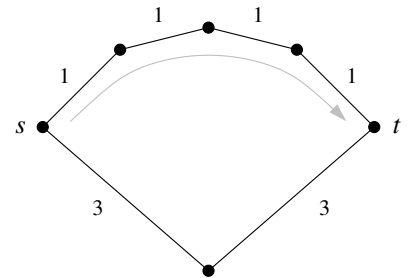


Figure 8.5: The shortest path in a network with varying edge lengths. The numbers on the edges in this network represent their lengths. The shortest path between s and t , taking the lengths into account, is the upper path marked with an arrow (which has total length 4), even though it traverses more edges than the alternative, lower path (which has length 6).

not be the same as the shortest path in terms of number of edges, as considered in Section 8.5.5. Take a look at Fig. 8.5, for example. The shortest path between s and t in this small network traverses four edges, but is still shorter, in terms of total edge length, than the competing path with just two edges. Thus we cannot find the shortest path in such a network using standard breadth-first search, which finds paths with the minimum number of edges. For problems like this we need a different algorithm. We need *Dijkstra's algorithm*.

Dijkstra's algorithm,¹⁷ like breadth-first search, finds the shortest distance from a given starting node s to every other node in the same component, but does so taking the lengths of edges into account.¹⁸ It works by keeping a record of the shortest distance it has found so far to each node and updating that record whenever a shorter one is found. It can be shown that, at the end of the algorithm, the shortest distance found to each node is in fact the shortest distance possible by any route. In detail the algorithm is as follows.

We start by creating an array of n elements to hold our current estimates of the distances from s to every node. At all times during the running of the algorithm these estimates are upper bounds on the true shortest distances, meaning the true distance to a node is less than or equal to the estimate. Initially we set our estimate of the distance from s to itself to be zero, which is trivially correct, and from s to every other node to be infinity, which is clearly a safe upper bound.

We also create another array of n elements in which we record when we are certain that the distance we have to a given node is in fact the exact shortest distance (and not merely a bound). For instance, we might use an integer array with 1s to indicate distances we are certain of and 0s for distances that are just our best current estimate. Initially, we put a 0 in every element of this array. (You might argue that we know for certain that the distance from s to itself is zero and hence that we should put a 1 in the element corresponding to node s .

¹⁷Named after its inventor, the Dutch computer scientist Edsger Dijkstra.

¹⁸We assume that the lengths are all strictly positive. If lengths can be negative, which happens in some networks, then the problem is much harder, falling in the class of "NP-complete" computational problems, for which even the best known algorithms take an amount of time exponential in n to finish, in the worst case [9]. Indeed, if edges are allowed to have negative lengths, there may not be any shortest path between a pair of nodes at all, since one can have a loop in the network that has negative length, so that one can reduce the length of a path arbitrarily by going around the loop repeatedly. We also assume there are no edges of length zero. We can define a meaningful shortest distance between nodes if there are length-zero edges, but the shortest *path* may be ill-defined, since there could be a loop of length zero that we could traverse an arbitrary number of times without adding to the length. To avoid all of these pathologies, therefore, we limit ourselves to strictly positive edge lengths.

Let us, however, pretend that we don't know this to begin with, as it makes the algorithm work out more neatly.)

Now we do the following:

1. Among all nodes about whose distance we are not yet certain, we find the node v that has the smallest estimated distance.
2. We mark this distance as being certain.
3. We calculate the distances from s via v to each of the neighbors of v by adding to v 's distance the lengths of the edges connecting v to each neighbor. If any of the resulting distances is smaller than the current estimated distance to the same node, the new distance replaces the older one.
4. We repeat from step 1 until the distances to all nodes are flagged as being certain.

Simple though it is to describe, it's not immediately obvious that this algorithm does what it is supposed to do and finds true shortest paths. The crucial step is step 2 where we declare the current smallest estimated distance in fact to be certain. That is, we claim that among nodes for which we don't yet definitely know the distance, the smallest distance recorded to any node is in fact the smallest possible distance to that node.

To see why this is true consider such a node, which we'll again call v . If v 's current estimated distance were not the true shortest distance, then there must exist some other path from s to v that has a shorter length. The situation is illustrated in Fig. 8.6.

Somewhere along this hypothetical shorter path there must exist a pair of adjacent nodes x, y such that x 's distance is known for certain and y 's is not. The current estimated distance for node y is at most equal to the shortest distance from s to x plus the length of the edge from x to y , since this estimate is calculated when we explore the neighbors of x in step 3 above. And this estimate is itself no greater than the distance from s to x to y along our hypothetical path, which in turn is no greater than the total length of the hypothetical path, which is strictly less than v 's estimated distance. In a notation in which e_{uv} is the current estimated distance from u to v and h_{uv} is the distance along the hypothetical path, we have:

$$e_{sy} \leq h_{sy} \leq h_{sv} < e_{sv}. \quad (8.5)$$

Thus, y 's estimated distance from s must be strictly less than v 's and we have a contradiction, since v is by hypothesis the node with the shortest estimated

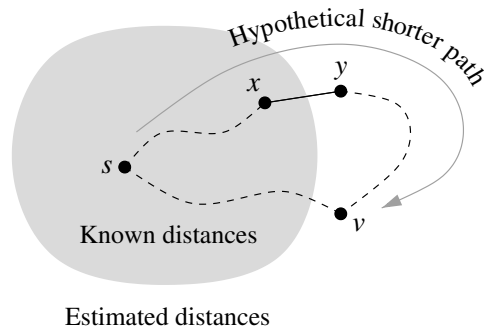


Figure 8.6: Proof of the correctness of Dijkstra's algorithm. If v is the node with the smallest estimated (i.e., not certain) distance from s , then that estimated distance must in fact be the true shortest distance to v . If it were not and there were a shorter path s, \dots, x, y, \dots, v then all points along that path must have shorter distances from s than v 's estimated distance, which means that y has a smaller estimated distance than v , which is impossible.

distance.¹⁹ Hence, there is no path to node v with length less than v 's current estimated distance and we can safely mark that distance as being certain, as in step 2 above.

Thus, on each round the algorithm correctly flags one additional distance as being known exactly and when all distances have been so flagged the algorithm ends.

As with breadth-first search, the running time of Dijkstra's algorithm depends on how it is implemented. The simplest implementation is one that searches through all nodes on each round of the algorithm to find the one with the smallest estimated distance. This search takes time $O(n)$. Then we must calculate a new estimated distance to each of the neighbors of the node we find, of which there are $O(m/n)$ on average. To leading order, one round thus takes time $O(m/n + n)$ and the whole algorithm, which runs (in the worst case) for n rounds, takes time $O(m + n^2)$ to find the distance from s to every other node.

But we can do better than this. The bottleneck in the algorithm is the process of searching for the node with the smallest estimated distance. We can speed up the calculation by storing the estimated distances in such a way that we always keep track of the smallest one, so we can find it quickly. To do this we

¹⁹One might imagine one could get around this if y and v were actually the same node, but this cannot be the case since we have just shown that y and v have different estimated distances, so they cannot be the same.

make use of a specialized data structure called a *heap*. A heap is an object that stores numbers—in this case the estimated distances—in such a way that the smallest number is always the first element in the heap. We will not describe here how a heap works, but the interested reader can find a description in many computer science texts, such as, for instance, Cormen *et al.* [122].

To implement Dijkstra’s algorithm using a heap, we start off by putting the estimated distances for all nodes into the heap (i.e., a single zero and $n - 1$ infinities). Then we repeatedly find and remove the node with the smallest estimated distance from the heap, explore its neighbors, and, if necessary, update their estimated distances, following the prescription given above. When all nodes have been removed from the heap and the heap is empty, the algorithm ends.²⁰ Because the smallest distance is always the first element in a heap we can find it in time $O(1)$, which is far faster than the $O(n)$ of the naive implementation. The price we pay for this, however, is that *removing* this node from the heap takes a slightly longer time $O(\log n)$. The operation of updating an estimated distance in the heap with a new and better estimate (which in the worst case we do an average of $O(m/n)$ times per round) also takes $O(\log n)$ time, and hence a complete round of the algorithm now takes time $O((m/n) \log n + \log n)$ and all n rounds take $O((m + n) \log n)$, or $O(n \log n)$ on a sparse network with $m \propto n$. This is very nearly the best running time known for this problem,²¹ and close to, though not quite as good as, the $O(m + n)$ for the equivalent problem on an unweighted network (factors of $\log n$ being close to constant given that the logarithm is a very slowly growing function of its argument).

As we have described it, Dijkstra’s algorithm finds shortest distances from node s to other nodes but, like breadth-first search, it can be modified also to find the actual paths that realize those distances. The modification is very similar to that for breadth-first search: we construct a shortest-path tree of pointers from nodes to their predecessors. We create such a pointer when we first assign a node an estimated distance less than infinity and move the pointer to point to a new node every time we find a new estimated distance that is less

²⁰If the network has more than one component then the algorithm ends when the smallest distance in the heap is infinity, since this tells us that the only distances left in the heap are to nodes in different components from s .

²¹There is more than one type of heap. The results described here are for the most common type, the *binary heap*. In theory one can achieve a slightly better running time of $O(m + n \log n)$ using another type of heap known as a *Fibonacci heap* [122], but in practice the operation of the Fibonacci heap is rather complicated and the calculation usually ends up running slower. Moreover, most programming languages already have binary heaps built in, in the form of library functions one can call, but few have Fibonacci heaps, so in practice Dijkstra’s algorithm is almost always implemented using a binary heap.

than the current one. The last position in which the pointer comes to rest indicates the true predecessor in the shortest-path tree. If a new estimate of the distance to a node is ever exactly the same as the current estimate then we add an additional pointer to the node, so that we have two pointers indicating the two alternative predecessors. When the algorithm is finished, the shortest-path tree, like those in Fig. 8.3 (page 250), can be used to reconstruct the shortest paths themselves, or to calculate other quantities such as a weighted version of betweenness centrality.

Dijkstra's algorithm finds a number of technological uses. For instance, it can be used to find the best routes for transmission of data over the Internet. Edge lengths in this case might represent travel time of data, say in microseconds, so that the shortest path is the one with shortest travel time. Mapping and navigation software also uses Dijkstra's algorithm to calculate the quickest driving or walking route to a given destination. Your phone or GPS device does this when you ask it for directions to a destination.

8.7 MAXIMUM FLOWS AND MINIMUM CUTS

In Section 6.13 we discussed independent paths, connectivity, cut sets, and maximum flows in networks. In particular, we defined two paths that connect the same nodes s and t to be edge-independent if they share none of the same edges and node-independent if they share none of the same nodes except for s and t themselves. The edge or node connectivity of the nodes is then the number of edge- or node-independent paths between them. We also showed that the edge or node connectivity is equal to the size of the minimum edge or node cut set—the minimum number of edges or nodes that need to be removed from the network to disconnect s from t . Connectivity is thus a simple measure of the robustness of the connection between a pair of nodes. Finally, we showed that the edge-connectivity is also equal to the maximum flow that can pass from s to t if we think of the network as a network of pipes, each of which can carry one unit of flow.

Thus, if we are able to find any one of three things—the number of independent paths between two nodes, the size of the minimum cut set, or the maximum flow—then we know all three, because they are all equal. It turns out that the easiest one to calculate is the maximum flow, and in practice all algorithms for solving such problems are in fact maximum flow algorithms. In this section we look at the most widely used maximum flow algorithm, the Ford–Fulkerson or augmenting path algorithm, which calculates the flow between two nodes in average time $O((m+n)m/n)$, and hence also calculates the number of edge-independent paths and the size of the minimum edge

cut set. With small extensions, the algorithm can also find the independent paths themselves or the specific set of edges that constitute the minimum cut set. A further simple modification of the algorithm allows us also to calculate node-independent paths and node cut sets.

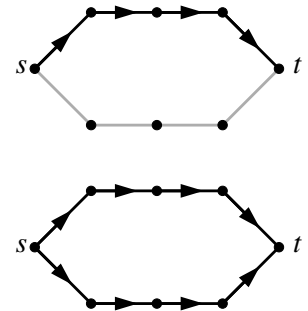
All the developments of this section are described for undirected networks, but the algorithms work perfectly well, without modification, for directed networks as well. Readers who want to know more about maximum flow algorithms are recommended to look at the book by Ahuja *et al.* [9], which contains hundreds of pages on the topic and covers almost every conceivable detail.

8.7.1 THE AUGMENTING PATH ALGORITHM

In this section we describe the augmenting path algorithm of Ford and Fulkerson for calculating maximum flows between nodes in a network.²² The case of primary interest to us is the one where each edge in the network can carry the same single unit of flow. The algorithm can be used in the more general case where the edges can have different capacities, but we will not discuss that case here.²³

The basic idea behind the augmenting path algorithm is a simple one. We first find a path from source s to target t using the breadth-first search algorithm of Section 8.5 and we imagine sending one unit of flow along this path.²⁴ This “uses up” some of the edges in the network, filling them to capacity so that they can carry no more flow. Then we find another path from s to t among the remaining edges and send another unit of flow along that. We repeat this procedure until no more paths can be found.

Unfortunately, this does not yet give us a working algorithm, because as we have described it the procedure will not always find the maximum flow. Consider Fig. 8.7a. If we apply breadth-first search between s and t we find the path marked in bold. Unfortunately, once we have used up all the edges



A simple breadth-first search finds a path from source s to target t (top) in this network. A second search using only the edges not used in the first finds a second path (bottom).

²²The augmenting path algorithm is not the only algorithm for calculating maximum flows. It is, however, the simplest and its average performance is about as good as any other, so it is a good choice for everyday calculations. It is worth noting, however, that the *worst-case* performance of the algorithm is quite poor—if one is particularly unlucky the algorithm can take a very long time to run. Another algorithm, the *preflow-push algorithm* [9], has much better worst-case performance and comparable average-case performance, but is considerably more complicated to implement.

²³See Ahuja *et al.* [9] or Cormen *et al.* [122] for details of the general case.

²⁴Technically, the augmenting path algorithm doesn’t specify how paths are to be found. Here we assume they are found using breadth-first search, which is known to give good performance. Sometimes this version of the algorithm is called the *shortest augmenting path algorithm* or the *Edmonds–Karp algorithm*.

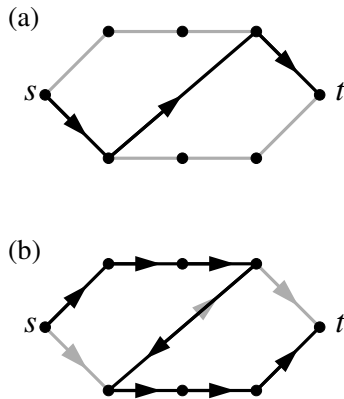


Figure 8.7: The augmenting path algorithm. (a) We find a first path from source s to target t using breadth-first search. This leaves no more independent paths from s to t among the remaining edges. (b) However, if we allow flows in both directions along an edge (such as the central edge in this network), then we can find another path.

along this path there are no more paths from s to t that can be constructed with the remaining edges, so the algorithm stops after finding just one path. It is clear, however, that there are in fact two edge-independent paths from s to t —along the top and bottom of the network—and a maximum flow of two, so the algorithm has given the wrong answer.

There is, however, a simple fix for this problem, which is to allow fluid to flow simultaneously *both ways* down an edge in our network. That is, we allow one unit of flow in each direction. If the edges were real pipes, this would not be possible: if a pipe is full of fluid flowing one way then there is no room for any to flow the other way. However, if fluid *were* flowing both ways down an edge, the net flow in and out of either end of that edge would be zero—the two flows would effectively cancel out, giving zero net flow. And zero flow down an edge certainly is possible.

So we use a trick and allow our algorithm to place a unit of flow both ways down any edge, but declare this to mean in practice that there is no flow at all on that edge. This means that the paths we find will no longer necessarily be independent paths, since two of them can share an edge so long as they pass along it in opposite directions. But this doesn't matter: the total amount of flow we find is still an allowed amount, since no pipe is ever required to carry more than one unit of flow. The paths found by the algorithm are called *augmenting paths*, to distinguish them from independent paths.

More generally, we can have any number of units flowing either way down an edge, provided they cancel out to give an allowed amount of net flow. Allowed net flows are (a) zero or (b) one unit of net flow in either direction. Thus, two units of flow one way down a pipe and two units the other way would be allowed, or three units one way and four the other, and so forth.

Three units one way and five the other would not be allowed, however.²⁵

To see how this works in practice, consider Fig. 8.7 again. We begin by performing a breadth-first search that finds the path shown in panel (a). Now, however, there is a second path to be found, as shown in panel (b), making use of the fact that we are still allowed to send one unit of flow *backwards* along the edge in the center of the network. After this, however, there are no more paths left from s to t , even allowing flows in both directions along edges, and so the algorithm stops and tells us that the maximum possible flow is two units, which is the correct answer.

This is merely one example of the algorithm: we still have to prove that it gives the correct answer in all cases, which we do in Section 8.7.3. First, however, let us look at how the algorithm is implemented and at its running time.

8.7.2 IMPLEMENTATION AND RUNNING TIME

Implementation of the augmenting path algorithm is straightforward. It just requires us to keep track of the amount and direction of flow along each edge. We begin by setting the flow along each edge to zero, then we carry out the following steps:

1. Perform a breadth-first search starting from node s and construct a shortest-path tree (Section 8.5.5). However, the breadth-first search is performed with the important constraint that when we examine the neighbors of a node we look only at those that are reachable along edges not yet filled to capacity in the direction we want to go. If the net flow along an edge is already one unit in the direction we are going then we do not follow that edge.
2. If node t is never reached by the breadth-first search, the algorithm ends.
3. Find a path between s and t in the shortest-path tree. If there is more than one such path then choose any of them—it doesn't matter which one.
4. Add one unit of net flow in the forward direction along each edge in this path.
5. Repeat from step 1.

At the end of the process the number of paths found from s to t is equal to the maximum flow.

²⁵On networks with directed edges, we allow either the same flow in both directions along an edge (i.e., zero net flow) or one more unit in the forward direction than in the backward direction, but not vice versa.

Each breadth-first search, including construction of the shortest-path tree and finding a path between s and t , takes time $O(m + n)$ for a network stored in adjacency list format (see Sections 8.3.2 and 8.5). The updates to the flows along the path take $O(m)$ time in the worst case where we have to update every edge in the network, so they do not change the overall $O(m + n)$ running time. As discussed in Section 6.13, the number of independent paths from node s to node t can be no greater than the smaller of the degrees k_s and k_t of the two nodes (since each path must leave or enter those nodes along a different edge). Thus the running time of the algorithm is $O(\min(k_s, k_t)(m + n))$. If we are interested in the average running time over many randomly chosen pairs of nodes, then we can make use of the fact that $\langle \min(k_s, k_t) \rangle \leq \langle k \rangle$ (where $\langle \dots \rangle$ denotes the average), and recalling that $\langle k \rangle = 2m/n$ (Eq. (6.15)), this implies that the average running time of the algorithm is $O((m + n)m/n)$, which is $O(n)$ on a sparse network with $m \propto n$. (On the other hand, on a dense network where $m \propto n^2$, we would have $O(n^3)$, which is much worse.)

In some cases, it is possible to improve the running time of the algorithm by using a two-source breadth-first search of the kind described in Section 8.5.4 instead of the one-source version described here. The worst-case running time of the two-source version is no better than the one-source version, but for specific networks, and particularly for the common case of networks with small diameter, it can give a boost in performance (see the discussion in Section 8.5.4).

8.7.3 WHY THE ALGORITHM GIVES CORRECT ANSWERS

It is plausible but not immediately obvious that the augmenting path algorithm correctly finds maximum flows. We can prove that it does as follows.

The augmenting path algorithm works by repeatedly finding augmenting paths, each of which contributes one unit of flow from s to t , so that the total flow found is equal to the number of paths. The algorithm stops when there are no more paths to be found. To prove that the algorithm is correct it suffices to show that the algorithm does not stop until the flow reaches its maximum. To put that another way, we need to prove the following statement:

If at some point in our algorithm the flow from s to t found so far is less than the maximum possible flow, then there must exist at least one more augmenting path in the network.

Consider such a point in the operation of the algorithm and let us represent the flows at that point by f , the set of all individual net flows along the edges of the network. And consider also the maximum possible flow from s to t , represented by f_{\max} , the corresponding set of individual net flows. (There may

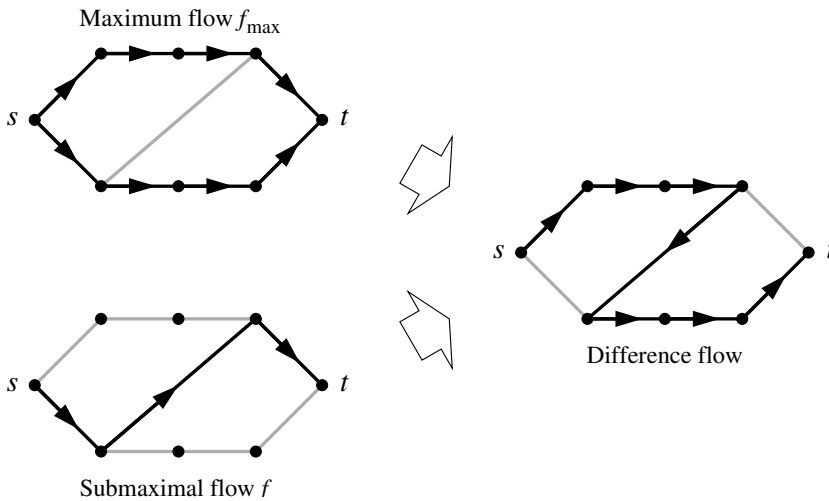


Figure 8.8: Correctness of the augmenting path algorithm. If we subtract from the maximum flow f_{\max} (upper left) any submaximal flow f (lower left) the resulting difference flow (right) necessarily contains at least one path from s to t , and that path is necessarily an augmenting path for f .

be more than one way of achieving the maximum flow, in which case f_{\max} can be any one of them we like—it doesn't matter which one.) By hypothesis, the total flow out of s and into t is greater in f_{\max} than in f . Let us calculate the difference flow $\Delta f = f_{\max} - f$, meaning we subtract the net flow along each edge in f from the net flow along the same edge in f_{\max} , respecting flow direction—see Fig. 8.8. (For instance, the difference of two unit flows in the same direction would be zero while the difference of two in opposite directions would be two in one direction or the other.)

Since the total flow is greater in f_{\max} than in f , the difference flow Δf must have a net flow out of s and a net flow into t . What's more, because the fluid composing the flow is conserved at nodes, every node except s and t must have zero net flow in or out in both f_{\max} and f and hence also in Δf , meaning that fluid is also conserved in Δf . But in that case the flow from s to t in Δf must form at least one path p across the network: if fluid leaves s , arrives at t , and is conserved everywhere in between, it must take *some* path across the network to get from s to t .

But if there is indeed such a path p , meaning that there is a flow in Δf in the forward direction along each edge in p , then there must have been no such flow in f along any of the same edges. If there were such a flow in f , then when we

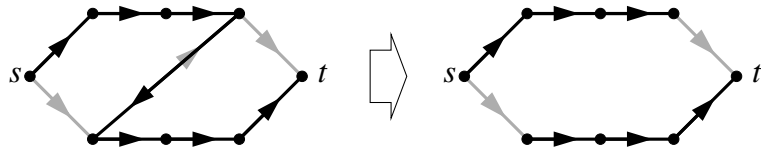


Figure 8.9: Reconstructing the independent paths from the final flows. Deleting every edge in the network that has zero net flow leaves a network consisting of the independent paths only.

performed the subtraction $\Delta f = f_{\max} - f$ the flow in Δf would be either zero or negative on the edge in question (depending on the flow in f_{\max}), but could not be positive. Thus we can always safely add to f a unit of flow forward along each edge in p without overloading any of the edges. This immediately implies that p is an augmenting path for f .

Thus, for any flow that is not maximal, at least one augmenting path always exists, and hence it follows that the augmenting path algorithm is correct and will always find the maximum flow.

8.7.4 FINDING INDEPENDENT PATHS AND MINIMUM CUT SETS

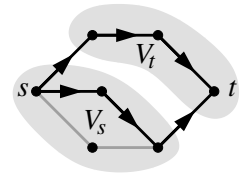
Once we have found the maximum possible flow between a given pair of nodes, we also automatically have the size of the minimum edge cut set and the number of edge-independent paths, which are both numerically equal to the number of units in the maximum flow (see Section 6.13).

We might also wish to know exactly where the independent paths run. The augmenting path algorithm does not give us this information directly since, as we have seen, the augmenting paths it finds are not necessarily independent paths. However, only a very small extension of the algorithm is necessary to find the independent paths: we take the final state of the network flows at the end of the algorithm and remove any edges from the network that have no net flow—see Fig. 8.9. In doing so we remove edges that have canceling flows in opposite directions and also edges that have no flow at all. The edges remaining are necessarily those that actually carry the maximum flow and it is a straightforward matter to trace these edges from s to t to reconstruct the paths taken by the flows.²⁶

²⁶Note that the independent paths are not necessarily unique: there can be more than one choice of paths, but this algorithm normally only finds one of them. Which one it finds depends on the particular choices made during the breadth-first search in situations where there is more

Another thing we might want is the set of edges that constitute the minimum cut set for the nodes s and t . In fact, in most cases there is more than one cut set of the minimum size, so more generally we would like to find one of the minimum cut sets. Again we can do this by a small extension of the augmenting path algorithm. We again consider the final pattern of flows in the network. By definition there is no way to add any additional flow from s to t to this pattern, since if there were the algorithm would not have stopped yet. If we try to perform an additional breadth-first search starting from s , we will in general be able to reach some nodes, but we won't be able to reach them all and, in particular, we will not reach t . Let V_s be the set of nodes reached by such a search and let V_t be the remaining nodes—those not in V_s , which necessarily includes node t (see figure). Then the set of edges that connect nodes in V_s to nodes in V_t constitutes a minimum cut set for s and t .²⁷

Why does this work? Clearly if we remove all edges that connect V_s and V_t we disconnect s and t , since then there is no path at all between them. Thus the edges between V_s and V_t do constitute a cut set. That it is a *minimum* cut set we can see by the following argument. Every edge from a node in V_s to a node in V_t must be carrying a unit of flow from V_s to V_t . If it were not, then it would have available capacity away from V_s , meaning that we would have followed it during our breadth-first search and the node at its other end would have been added to V_s , instead of being in V_t . But if every edge between V_s and V_t is carrying a unit of flow, the number of such edges—the size of the cut set—is equal to the size of the flow from V_s to V_t , which is also the flow from s to t . And, by the max-flow/min-cut theorem, a cut set between s and t that is equal in size to the maximum flow is a minimum cut set, and hence our result is proved.



The sets V_s and V_t for a small network.

8.7.5 NODE-INDEPENDENT PATHS

Once we know how to find edge-independent paths it is straightforward to find node-independent paths as well. First, note that any set of node-independent

than one path from s to t in the shortest-path tree. Furthermore, there can be points in the network where edge-independent paths come together at a node and then part ways again. If such points exist, we will have to make a choice about which way to go at the parting point. It doesn't matter what choice we make in the sense that all choices lead to a correct set of paths, but different choices will give different sets of paths.

²⁷A nice feature of this procedure is that normally one has to perform an additional breadth-first search anyway as part of the main augmenting path algorithm, since that is how one knows the algorithm is finished—it fails to find another augmenting path. So all the hard work has already been done by the main algorithm and finding the cut set is a minor extra step.

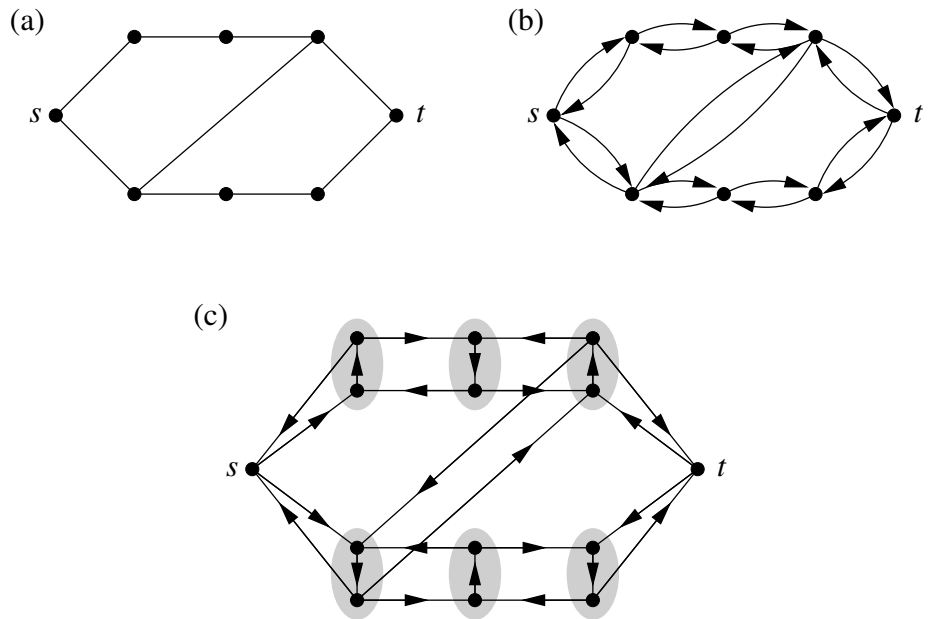


Figure 8.10: Mapping from the node-independent path problem to the edge-independent path problem. Starting with an undirected network (a), we (b) replace each edge by two directed edges, then (c) replace each node, except for s and t , with a pair of nodes with a directed edge between them (shaded) following the prescription in Fig. 8.11. Edge-independent paths on the final network then correspond to node-independent paths on the initial network.

paths between two nodes s and t is necessarily also a set of edge-independent paths: if two paths share none of the same nodes, then they also share none of the same edges. Thus, we can find node-independent paths using the same algorithm we used to find edge-independent paths if we just add the restriction that no two paths may pass through the same node, or equivalently that at most one unit of flow can pass through each node. One way to impose this restriction is the following. First, we replace our undirected network with a directed one, as shown in Fig. 8.10, with a directed edge in either direction between every connected pair of nodes. This does not change the maximum flow possible in the network and hence does not change the number of independent paths either.

Second, we replace each of the nodes in the network, except s and t , with a construct like that shown in Fig. 8.11. Each node is replaced with two nodes

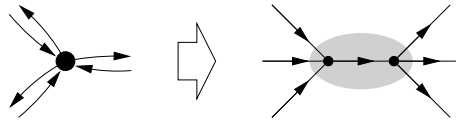


Figure 8.11: Node transformation for the node-independent path algorithm. Each node in the network is replaced by a pair of nodes joined by a single directed edge. All incoming edges are connected to one of the pair and all outgoing edges to the other as shown.

separated by a directed edge. All original incoming edges connect to the first of these two (on the left in Fig. 8.11) and all outgoing edges to the second. This new construct functions as the original node did, allowing flows to pass in along ingoing edges and out along outgoing ones, but with one important difference: assuming that the new edge joining the two nodes has unit capacity like all others, we are now limited to just one unit of flow through the entire construct, since every path through the construct must traverse this central edge. Thus, every allowed flow on this network corresponds to a flow on the original network with at most a single unit passing through each node.

Transforming the entire network of Fig. 8.10b using this method gives us a network that looks like Fig. 8.10c. Now we simply apply the normal augmenting path algorithm to this directed network and the number of *edge*-independent paths we find is equal to the number of *node*-independent paths on the original network of Fig. 8.10a.

In this chapter we have looked at a selection of the best known and most widely used network algorithms. There are many further algorithms that will be of interest to us in this book, but rather than collect all of them together here, we will instead introduce them in context as they arise during our continuing discussion of the theory and practice of networks. Armed with the fundamental tools and experience gained in this chapter, you should have no difficulty following the working of any of the algorithms that appear later on.

EXERCISES

8.1 What (roughly) is the time complexity of:

- a) Vacuuming a carpet if the size of the input to the operation is the number n of square feet of carpet?
- b) Finding a word in a (paper) dictionary if the size of the input is the number n of words in the dictionary?

8.2 Consider the following situations:

- a) You are asked to calculate the closeness centrality of a single node in an undirected network with m edges and n nodes. What algorithm would you use to do this, and what would be the time complexity of the operation in terms of m and n ?
- b) You are given a road map and told the average driving time along each road segment, then you are asked to find the route from A to B with the shortest average driving time. What algorithm would you use to do this, and what would be the time complexity of the calculation?
- c) What algorithm would you use to find all the components in an undirected network, and what would be the time complexity of the operation?
- d) What algorithm would you use to determine whether there are at least two node-independent paths between a given pair of nodes? Hence or otherwise, suggest an algorithm that can find all the bicomponents in a network.

8.3 Suppose you have a sparse undirected network with $m \propto n$. What is the time complexity of:

- a) Multiplying an arbitrary n -element vector by the adjacency matrix, if the network is stored in adjacency matrix format.
- b) The same multiplication if the network is in adjacency list format.
- c) The “modularity matrix” \mathbf{B} of an undirected network is the $n \times n$ symmetric matrix with elements

$$B_{ij} = A_{ij} - \frac{k_i k_j}{2m}.$$

(See Eq. (14.2) on page 500.) What is the time complexity for multiplying an arbitrary vector by the modularity matrix of our sparse network if the network is in adjacency list format? Describe briefly an algorithm that achieves this. (Hint: One’s first guess might be that the multiplication takes time $O(n^2)$ but it is possible to do it faster than this in the sparse network case.)

8.4 An interesting question, which is discussed in some detail in Chapter 15, concerns what happens to a network if you disable or remove its nodes one by one. The question is of relevance, for instance, to the vaccination of populations against the spread of disease. One typical approach is to remove nodes in order of their degrees, starting with the highest degrees first. Note that once you remove one node (along with its associated edges) the degrees of some of the other nodes may change.

In most cases it is not possible to do the experiment of removing nodes from a real network to see what effect it has, but we can simulate the process on a computer by

taking a network stored in computer memory, removing some of its nodes, and then measuring various properties of the remaining network.

- a) What is the time complexity for finding the highest-degree node in a network, assuming the nodes are given to you in no particular order?
- b) If we perform the repeated node removal in a naive way, searching exhaustively for the highest-degree node, removing it, then searching for the next highest, and so forth, what is the time complexity for removing all of the nodes? You can assume the network is stored in adjacency list format.
- c) Describe how the same operation could be performed with the degrees of the nodes stored instead in a heap. What now is the time complexity of the entire calculation?
- d) Taking the same approach, describe in a sentence or two a method for taking n numbers in random order and sorting them into decreasing order using a heap. Show that the time complexity of this sorting algorithm (which is called *heapsort*) is $O(n \log n)$.
- e) The degrees of the nodes in a simple network are integers between zero and n . It is possible to sort such a set of integers into decreasing (or increasing) order in time $O(n)$. Describe briefly an algorithm that achieves this feat.

8.5 What is the time complexity, as a function of the number n of nodes and m of edges, of the following network operations if the network in question is stored in adjacency list format?

- a) Calculating the mean degree.
- b) Calculating the median degree.
- c) Calculating the air-travel route between two airports that has the shortest total flying time, assuming the flying time of each individual flight is known.
- d) Calculating the minimum number of routers that would have to fail to disconnect two given nodes on the Internet.

8.6 For the breadth-first search algorithm of Section 8.5, prove the two statements given on page 242. That is, assuming we know the true shortest distance from a given starting node s in a network to all nodes at distance d or less, and that the distance to all other nodes has not yet been calculated, prove the following:

- a) For a node at distance d , every neighbor that has not yet been assigned a distance must have distance $d + 1$.
- b) Every node at distance $d + 1$ is a neighbor of at least one node with distance d .

Hint: While statement 1 should be straightforward to prove, statement 2 is a little more involved. It may be useful to consider a path from s to a node at distance $d + 1$, and then consider the shortest distance from s to the penultimate node along that path.

8.7 For an undirected network of n nodes and m edges stored in adjacency list format show that:

- a) It takes time $O(n(m + n))$ to find the diameter of the network.

b) It takes time $O(\langle k \rangle)$ on average to list the neighbors of a node, where $\langle k \rangle$ is the average degree in the network, but time $O(\langle k^2 \rangle)$ to list the second neighbors.

8.8 For a directed network in which in- and out-degrees are uncorrelated, show that it takes time $O(m^2/n)$ to calculate the reciprocity of the network. Why is the restriction to uncorrelated degrees necessary? What could happen if they were correlated?

8.9 Suppose that we define a new centrality measure x_i for node i in a network to be a sum of contributions as follows: 1 for node i itself, α for each node at (shortest) distance 1 from i , α^2 for each node at distance 2, and so forth, where $\alpha < 1$ is a given constant.

- a) Write an expression for x_i in terms of α and the shortest distances d_{ij} between node pairs.
- b) Describe briefly an algorithm for calculating this centrality measure.
- c) What is the time complexity of calculating x_i for all i ?

CHAPTER 9

NETWORK STATISTICS AND MEASUREMENT ERROR

*A discussion of the statistics of network measurements
and the types of errors that can arise in network data*

WHEN making empirical measurements of network structure, an important but frequently overlooked issue is the possibility that the measurements may not be completely accurate—they may contain errors. With most scientific experiments or observations, it is customary to report not only the measurements we make but also the estimated size of our errors. When measuring a voltage or a chemical concentration in the lab, we report the observed value plus a standard deviation or confidence interval. When performing a behavioral experiment or measuring the effectiveness of a medical treatment, we report the outcome along with a p -value or other statistic indicating our confidence in the result. In measurements of network structure of the kind described in Chapters 2 to 5, however, it is surprisingly common to omit any mention of the reliability of the results. Yet at the same time it is reasonable to suppose that there are errors in the measurements. Experimental data on the structure of biological networks, such as protein interaction networks or neural networks, for instance, are susceptible to measurement error in the lab [283, 459, 473]. Measurements of the structure of the Internet, using traceroute or BGP tables, suffer from incomplete sampling and technical limitations [108, 284]. Measurements of social network structure can be affected by subjectivity on the part of both participants and experimenters, recording error, and various kinds of

measurement error as well [56,259,320].

Sometimes the absence of error estimates on network data can be justified by defining the network in terms of the data. For instance, surveys of who people claim as their friends are an error-prone way of measuring actual friendship networks. But if one defines the object of interest to be the network of who *says* they are friends with whom, then by definition the data are a good representation. This, however, is a somewhat unsatisfactory approach. Most of the analyses we would want to perform on such a network implicitly assume that the data reflect actual friendships, not just reported ones, and to the extent that the two differ our analyses will be in error.

It has been shown that errors on network data can have a big impact on the accuracy of the conclusions drawn from those data [77,163,171,276,461]. There are many published network studies whose results might be called into question if one were to perform a careful analysis of the errors in the data. In this chapter we consider the various types of error that can occur in network data, along with statistical techniques for representing and quantifying them and understanding their impact on our analyses.

9.1 TYPES OF ERROR

Suppose we have measured a network of some kind, such as a technological, social, or biological network. In what ways could the resulting data be in error? Let us first consider the basic case of an undirected, unweighted simple network: each pair of nodes is either connected by a single undirected edge or not, and the network can be represented by a symmetric adjacency matrix \mathbf{A} with elements $A_{ij} = 0$ or 1.

Possible errors in such a network can be divided into errors on the nodes and errors on the edges. Potential node errors include the following:

Missing nodes: A common type of error is the omission of one or more nodes from a network. There might be a species in a food web that was not observed by the experimenter. There might be an individual omitted from a social network because he or she failed to complete a survey or questionnaire when asked. There might be a web page missing from a crawl of the Web because there were no links to it, so it was never found.

Erroneous extra nodes: A rarer error is a node that appears in a network but does not exist in reality. Examples of this type of error might be a web page that has been deleted but is erroneously still included in the network, or an individual in a social network who does not meet the criteria for inclusion (such as age or location).

Extra copies of nodes: A special case in which extra nodes are more common is when we mistakenly represent a single node by two or more different nodes. Such errors are particularly frequent in social network studies, especially when individuals are identified by name alone. In studies of academic coauthorship networks of the kind discussed in Section 4.5, for instance, authors of papers are normally identified by their published names. It is not uncommon for an author to use slightly different versions of their name on different publications—Frank Lloyd Wright, Frank Wright, F. L. Wright, and so forth—and this gives rise to errors when separate nodes are incorrectly created to represent each of these names even though they refer to the same person. The process of trying to determine when two different nodes actually represent the same person is called “entity resolution” or “node disambiguation,” and is discussed in Section 9.4.2.

Erroneously aggregated nodes: The flipside of the problem of extra copies of nodes is erroneous aggregation of nodes, where two different nodes are mistakenly combined into one. Taking the example of a coauthorship network again, it might happen that two different authors have the same name, in which case their publication records might get combined so that they appear as a single node in the network instead of two.

There can also be errors in the edges of a network:

Erroneously omitted edges: Edges that should be present may be omitted from the data. These are *false negatives*.

Erroneous extra edges: Edges may be reported as present that do not actually exist. These are *false positives*.

Missing data: We may lack any data on the presence or absence of an edge between two nodes. If a particular connection is simply never measured, then we do not know whether there is an edge there or not. Often this kind of omission is treated as a non-edge, but it is important to appreciate that the two are not the same: absence of evidence is not evidence of absence. By assuming that there is no edge between two nodes when in fact we simply have no information, there is a chance that we may be introducing a false negative.

Errors on edges are easier to treat statistically and better understood than errors on nodes, and for these reasons we will spend more time on them than on node errors. Node errors are, nonetheless, common and could potentially have a significant impact on the outcome of our studies and calculations, so we should not ignore their existence.

There are also additional types of errors that can occur in directed or weighted networks. The direction of an edge in a directed network could, for instance,

be misreported, although this is a relatively rare problem in most cases. A more common issue is the misreporting or misestimation of weights in weighted networks. Experimental error could easily affect measurements such as the volume of traffic on a road or the energy flow between species in a food web.

9.2 SOURCES OF ERROR

There are many different sources of error in measurements of network structure, depending on the the type of the network and the type of measurement. We won't try to list every possible error, but here are some illustrative examples.

Social networks measured using surveys or questionnaires: In social networks such as friendship networks whose structure is determined by the administration of surveys, interviews, or questionnaires, principal sources of error include subjectivity or bias on the part of the interviewer or interviewee, quantification and recording error, and missing data [56, 259, 320]. In asking who someone's friends are, for instance, different respondents may interpret the word "friend" in different ways and hence one person might classify a certain relationship as a friendship where another would not. Subjectivity on the part of the experimenter can arise, for example, when respondents give ambiguous answers and some interpretation is required to decide what those answers mean. Experimenters go to some length in the design of surveys to reduce potential subjectivity, but it is inevitable that some remains. Quantification errors can arise in the process of taking qualitative answers given by participants and turning them into hard numbers, while recording errors arise when, for one reason or another, the response to a question is recorded incorrectly.

Missing data is also a common problem in social surveys. Participation in such surveys is usually voluntary and participants may decline to answer some or all questions. Some members of the target group might be missing altogether: a survey of a company for instance, might miss some individuals who happen to be absent on the day the survey is conducted. Especially in so-called longitudinal studies, meaning those that measure the same network repeatedly over time, missing data can be a big problem. Getting the same group of participants to take repeated surveys over an extended period can be challenging, and there is often a significant fraction of individuals who are present for the first rounds but fail to turn up for later ones.

A more pernicious form of error in survey data, one that can be difficult to deal with, is the omission or truncation of responses as a result of specific design features in the survey itself. If a survey or questionnaire simply fails to ask about certain things, then those data will be missing. A classic example

arises in “fixed choice” studies of the type described in Section 4.2, in which respondents are asked to list a limited number of network contacts. We looked, for instance, at the study by Rapoport and Horvath [400] of friendship networks among schoolchildren, in which participants were asked to name up to eight of their best friends in school, but no more. Any participants who had more than eight friends had no opportunity to say so, and it is likely that many friendships went unreported as a result.

Social networks constructed from other data: Social networks constructed by direct observation, such as animal social networks, are less subject to bias on the part of the observed, who do not actively participate in the data gathering, but there can still be bias on the part of the experimenter in interpreting observations, as well as quantification and recording errors. Social networks constructed from archival or third-party records are quite reliable in some ways but suffer from their own problems. In the network of who is “friends” with whom in an online social network like Facebook, for instance, there is in a sense no ambiguity about which nodes are connected to which. Either two people are Facebook friends or they aren’t. However, different people may have different standards for making or accepting friendship requests. Some people may have a high threshold and accept only a few requests while others may be happy to be friends with everyone who asks. If the experimenter nonetheless assumes that all edges in the network represent equivalent relationships then conclusions drawn as a result may be in error. Data from third-party and archival sources can also be incomplete and can display erroneous aggregation or disaggregation of nodes of the kind described in Section 9.1. A network of who emails whom, for instance, is only as good as the email data available, which might be incomplete or cover only a limited span of time. It might also include a significant fraction of duplicate nodes if individuals are identified by their email addresses alone; it is not uncommon for one individual to have two or more email addresses.

Biological networks: Biological networks, such as the metabolic, protein, and regulatory networks of Section 5.1, are measured in the lab and their primary source of error is conventional laboratory measurement error. For a wide range of reasons, including natural variation in biological systems and inconsistent measurement conditions, experiments usually don’t give exactly the same results every time you perform them, meaning that any individual measurement may be in error, potentially by a large amount.

Consider, as an example, the protein–protein interaction networks discussed in Section 5.1.2. As described in that section, there are several methods for determining whether two proteins interact to form a protein complex, such as

co-immunoprecipitation, two-hybrid screens, and tandem affinity purification. These methods are of varying degrees of accuracy, but none is wholly reliable. In practice, this means that if we simply repeat the same experiment twice it may not give the same result. One time it may tell us that two proteins interact and another time that they do not. For example, Krogan *et al.* [283] assembled an extensive network of protein interactions using over 4000 separate affinity purification experiments, combining the data to give an overall picture of the network. In principle, one could imagine imposing a simple standard that says if a particular interaction is found to be present in, say, a half or more of the experiments that probe it then it will be considered correct and added to the network. In practice, however, there are more sophisticated and reliable ways to interpret data of this kind—see Section 9.3.

The Internet: Empirical studies of the structure of the Internet are susceptible to their own special types of measurement error. As discussed in Section 2.1, Internet structure is typically measured by finding paths between nodes on the network (either using traceroute or BGP tables) and then aggregating those paths to create a picture of the network. A problem with this approach is that it omits any edge that appears in none of the sampled paths. Even for very large sample sizes it is likely that some edges will be omitted and the problem is compounded by the fact that for practical reasons studies tend to sample many paths from the same starting point, rather than from different starting points. As argued by Lakhina *et al.* [284], the probability of a particular edge falling on a sampled path decreases with distance from the starting point, and hence for a set of paths starting from a single point, nearby edges can be significantly more likely to be discovered than those far away. Clauset, Moore, and co-workers [3, 108] have shown that this process can, for instance, generate the false impression that the network has a power-law degree distribution when in fact it does not, a finding that calls into question one of the most iconic results about the structure of the Internet, namely that its degree distribution approximately follows a power law (see Section 10.4).

Degree distributions, one of the most fundamental properties of networks, are discussed in Section 10.3.

The World Wide Web: Observations of hyperlinks on the World Wide Web are, in principle, highly reliable. If a link between two web pages is observed then there is unlikely to be any debate about it. On the other hand, links appear and disappear over time, so a measurement today may not be a good indication of the state of the network tomorrow or yesterday. Moreover, measurements suffer from the problems created by dynamic pages, as discussed in Section 3.1. Many of the pages we see on the Web are dynamically generated from databases and only come into existence when someone asks for them—the pages of search results returned by search engines like Google are a good example. Since the

number of possible dynamic pages is effectively infinite, one cannot realistically map them all. So in constructing maps of the web network one must make some decision about when and how to include dynamic pages, which introduces a level of subjectivity and arbitrariness into the network structure.

But perhaps the most serious source of error in measurements of the Web is the wholesale omission of inaccessible regions of the network. As discussed in Section 3.1, the web network is typically measured by automated surfing of web links to find pages, and this process by definition will not find pages that either (a) no one links to or (b) are linked to only by other pages that are not found. As a result any region of the network that is “upstream” of the current location cannot be reached by surfing—an in-component in network nomenclature—along with any region that is disconnected completely from the rest of the network. It is by definition difficult to estimate the size of the unobserved portion of the Web, but some studies put it at as much as a half of the network [84].

See Section 10.1.1 for a discussion of component structure and the reachability of sites on the Web.

9.3 ESTIMATING ERRORS

Let us turn now to methods for quantifying and estimating errors and the effects they can have on the network measurements we care about. It will be useful to start our discussion by first revisiting the techniques used to quantify ordinary measurement errors on real quantities, such as lengths, weights, voltages, and so forth.

9.3.1 CONVENTIONAL STATISTICS OF MEASUREMENT ERROR

When we measure some quantity x in the lab or in the field we implicitly assume that there is some true underlying value of that quantity, sometimes called the *ground-truth* value, which is unknown. Let us denote this value z . Our measurements are usually an imperfect reflection of this ground-truth value—there is some process, which we often do not understand completely, that introduces error into our results. The standard way to deal with this situation is to assume a model, often somewhat simplified, to describe this process. The most common model by far is the Gaussian or normal model, in which the measured value x is assumed to be equal to the true value z plus a random additional amount—the error—drawn from a normal distribution.

Mathematically, we say that the measured value x is drawn from the probability distribution

$$P(x|z, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-z)^2/2\sigma^2}, \quad (9.1)$$

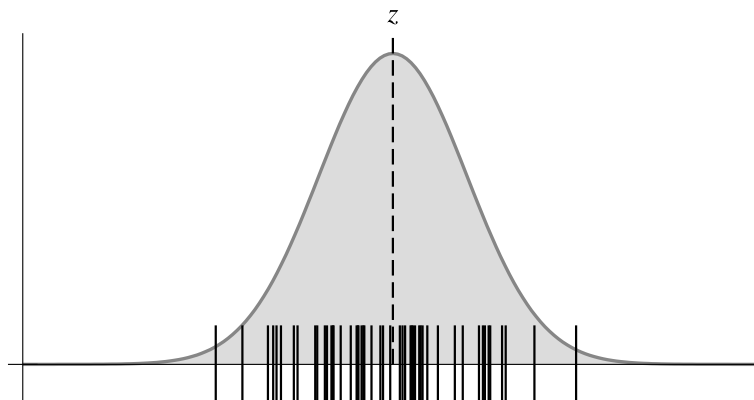


Figure 9.1: Measurement error on a real-valued quantity. We repeatedly measure some quantity whose true value z falls at the dashed line. But our measurements, represented by the lines on the horizontal axis, contain experimental error and hence are not exactly equal to z , or to each other. Instead, they are distributed randomly around z following, in this case, a normal or Gaussian distribution (solid curve). Our goal is to estimate the value of z and the width of the curve from the data.

where the initial factor of $1/\sqrt{2\pi\sigma^2}$ is a normalizing constant that ensures that the distribution integrates to unity, and σ is the standard deviation of the normal distribution. The value of σ parametrizes the size of the error added to the ground-truth value.

The goal of our experiments is to make the best estimate we can of the ground-truth value z , along with an estimate of σ , which tells us the typical size of the error. The standard way to do this is to make several measurements of x and then combine them to estimate z and σ . The theory behind how we do this is based on the *method of maximum likelihood*.

9.3.2 THE METHOD OF MAXIMUM LIKELIHOOD

Suppose we make N measurements of a quantity of interest and let us denote the results of these measurements by x_1, \dots, x_N . Our assumption will be that these measurements are (a) statistically independent, meaning that the value of one measurement has no effect on any other, and (b) distributed randomly around the true value z following a normal distribution with standard deviation σ , as in Eq. (9.1). The situation is depicted in Fig. 9.1.

Since the measurements are independent, the probability of all of them,

collectively, is simply the product of their individual probabilities:¹

$$P(x_1, \dots, x_N | z, \sigma) = \prod_{i=1}^N P(x_i | z, \sigma) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - z)^2 / 2\sigma^2}. \quad (9.2)$$

This probability is called the *likelihood* of the data x_i given z and σ .

Normally, we do not know the values of z and σ . We have only the observed data x_i . Can we nonetheless estimate the values of z and σ ? Indeed we can. Specifically, we can ask what values of z and σ are most likely given the data. The probability that particular values z and σ were responsible for generating the x_i can be calculated using Bayes' rule for probabilities thus:

$$P(z, \sigma | x_1, \dots, x_N) = P(x_1, \dots, x_N | z, \sigma) \frac{P(z)P(\sigma)}{P(x_1, \dots, x_N)}, \quad (9.3)$$

where $P(z)$ is the so-called *prior probability* of z , often just called “the prior” for short. It is the probability of a particular value of z if we don't know (or have not yet measured) the data x_i . Similarly, $P(\sigma)$ and $P(x_1, \dots, x_N)$ are the prior probabilities of σ and the data.

The most likely values of z and σ are, by definition, those with the highest probability $P(z, \sigma | x_1, \dots, x_N)$, which can be found by maximizing Eq. (9.3) with x_1, \dots, x_N held fixed at the observed values. But if the x_i are held fixed, then the prior $P(x_1, \dots, x_N)$ is fixed too—it is just a constant. Moreover, one commonly assumes that the priors $P(z)$ and $P(\sigma)$ are also constant, i.e., that all values of these quantities are a priori equally likely.² With these assumptions we have

$$P(z, \sigma | x_1, \dots, x_N) \propto P(x_1, \dots, x_N | z, \sigma). \quad (9.4)$$

But this implies that the maximum of $P(z, \sigma | x_1, \dots, x_N)$ is in the same place as the maximum of $P(x_1, \dots, x_N | z, \sigma)$. Hence the best values of z and σ are obtained simply by maximizing the likelihood, Eq. (9.2), with respect to both parameters.

This is the method of maximum likelihood. Applied to the current problem it involves simply maximizing Eq. (9.2) in standard fashion by differentiating and setting the result to zero. For instance, differentiating with respect to z

¹Technically, this is a *probability density*—the probability of a value in the small element of volume $d^N \mathbf{x}$ is $P(x_1, \dots, x_N | z, \sigma) d^N \mathbf{x}$ —but the argument is the same either way.

²We cannot really assume that all values are equally likely, since both z and σ have infinite range. But we could assume that all values are equally likely over some very large range, encompassing everything we can reasonably expect to encounter with the particular data we are considering.

gives

$$\sum_{i=1}^N \frac{x_i - z}{\sigma^2} \prod_{j=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_j - z)^2 / 2\sigma^2} = 0. \quad (9.5)$$

Canceling several factors, we can simplify this to $\sum_i (x_i - z) = 0$, or equivalently

$$z = \frac{1}{N} \sum_{i=1}^N x_i. \quad (9.6)$$

Similarly, if we differentiate with respect to σ we get

$$\sum_{i=1}^N \left[\frac{(x_i - z)^2}{\sigma^3} - \frac{1}{\sigma} \right] \prod_{j=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_j - z)^2 / 2\sigma^2} = 0, \quad (9.7)$$

which simplifies to $\sum_i [(x_i - z)^2 - \sigma^2] = 0$ or

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - z)^2. \quad (9.8)$$

In other words, the best estimates of the mean z and standard deviation σ of the Gaussian distribution, based on the data, are just given by the familiar standard formulas, Eqs. (9.6) and (9.8).

All of us learned these formulas as students, but no one ever says why these are the right formulas to use. They just expect you to take them on faith. But you don't have to: as we see here, the formulas for the mean and standard deviation can be derived using the principle of maximum likelihood.

There is one small further wrinkle on the method that will be useful to us. In most practical situations, we work not with the likelihood itself but with its logarithm, often called the *log-likelihood* for short. In the present case, the log-likelihood is

$$\log P(x_1, \dots, x_N | z, \sigma) = -\frac{1}{2} N \log 2\pi\sigma^2 - \sum_{i=1}^N \frac{(x_i - z)^2}{2\sigma^2}. \quad (9.9)$$

Since $\log x$ is a monotone increasing function of its argument—bigger x always means bigger $\log x$ —the maximum of $\log x$ falls in the same place as the maximum of x . So we will get the same answer for the position of the maximum whether we maximize the likelihood or the log-likelihood. In practice maximizing the log-likelihood is almost always easier. In the present case, for

instance, instead of the complicated expression in Eq. (9.5), differentiating (9.9) with respect to z and setting the result to zero gives simply

$$\frac{1}{\sigma^2} \sum_{i=1}^N (x_i - z) = 0, \quad (9.10)$$

which trivially rearranges to give Eq. (9.6) again, and a similar calculation works for σ too.

Another simplifying observation is that constants that multiply the likelihood do not change the position of its maximum, and hence can be ignored. Such constants correspond to additive terms in the log-likelihood. For instance, the first term in Eq. (9.9) can be rewritten as

$$\frac{1}{2}N \log 2\pi\sigma^2 = \frac{1}{2}N \log 2\pi + N \log \sigma. \quad (9.11)$$

The quantity $\frac{1}{2}N \log 2\pi$ is constant and hence can be ignored. To put that another way, when we differentiate with respect to z or σ to find the maximum, this term disappears and hence it plays no role in our calculation.

9.3.3 ERRORS IN NETWORK DATA

Having seen how errors are treated in ordinary experimental data, let us now return to the issue of errors in network data. We will apply the same principles, assuming a set of measurements and a model for the error and then employing the method of maximum likelihood to estimate both the ground truth and the error. Our discussion closely follows that of Refs. [361, 362].

When we measure a network the measured object is more complicated than a single number; it is the entire adjacency matrix of the network. Nonetheless, we can imagine describing our uncertainty about our measurements in the same way using a probabilistic model. We assume that there is some underlying ground-truth network represented by an adjacency matrix \mathbf{A} , which is unknown. Then we make some measurements of the network structure, which potentially contain measurement error of some kind, and we would like to estimate \mathbf{A} or an approximation to it from these measurements. There are many forms measurements could take, including measurements of the entire structure of the network, measurements of individual edges, or measurements of paths such as metabolic pathways or traceroute paths on the Internet. In addition we may perform repeated measurements of some or all elements, to get a handle on how much they vary as a result of error.

We also assume a model that describes how error is introduced into our measurements. By contrast with the case of simple real-valued data studied in

the preceding section, there is no one standard model that is broadly applicable in most cases. Models for network error take different forms, depending on the type of network we are considering and the kind of measurements we are making. We discuss some specific examples of error models in the following sections, but for the moment let us keep our discussion general. The model specifies the probability of making a particular set of measurements—the data—given the ground truth \mathbf{A} and optionally some additional parameters, analogous to the standard deviation σ in the case of real-valued data. Let us denote this probability $P(\text{data}|\mathbf{A}, \theta)$, where θ indicates all of the parameters.

Now, as previously, we employ Bayes' rule to write

$$P(\mathbf{A}, \theta|\text{data}) = P(\text{data}|\mathbf{A}, \theta) \frac{P(\mathbf{A})P(\theta)}{P(\text{data})}. \quad (9.12)$$

Assuming we can write down expressions for all the quantities on the right-hand side, we can now maximize this probability to find the most likely values of the matrix \mathbf{A} and the model parameters θ . These, along with the model itself, give us not only an estimate of the structure of the network but also a quantification of the error on that structure.

9.3.4 THE EM ALGORITHM

If we were to take an approach directly analogous to that of Section 9.3.2 we would simultaneously maximize (9.12) with respect to both \mathbf{A} and θ . Even for simple error models, however, this calculation can be difficult. The matrix \mathbf{A} is a discrete-valued object, so we cannot maximize by simple differentiation as we did for our Gaussian error model, and while the parameters θ might be continuous-valued, derivatives with respect to them often result in complicated equations that are hard to solve.

Instead, therefore, we take a different approach that leads to one of the most widely used and elegant methods in statistics, the *expectation-maximization algorithm*, also called the EM algorithm for short.

As the first step in our calculation, we will try to find a value for the parameter or parameters θ , neglecting \mathbf{A} for the moment. For simplicity, let us assume that there is only one parameter θ . The generalization to the case of two or more is straightforward.

We can find the best estimate of θ by first writing

$$P(\theta|\text{data}) = \sum_{\mathbf{A}} P(\mathbf{A}, \theta|\text{data}). \quad (9.13)$$

Here the sum is over all possible adjacency matrices \mathbf{A} , i.e., over all networks with the same number n of nodes as the observed network. (We will continue

to focus on simple undirected unweighted networks, although in principle the calculation can be generalized to other cases.) Maximizing $P(\theta|\text{data})$ over θ now gives us the most probable value of θ given the observed data. As discussed in Section 9.3.2, we commonly actually maximize the logarithm of the probability, not the probability itself, and we will do that here too. That is, we will maximize

$$\log P(\theta|\text{data}) = \log \sum_{\mathbf{A}} P(\mathbf{A}, \theta|\text{data}) \quad (9.14)$$

with respect to θ .

As we have said, direct maximization turns out to be difficult in this case—the expressions we get are complicated and hard to work with. Instead, therefore, we employ the following trick. We make use of *Jensen's inequality*, which says that the log of the weighted average of a set of quantities is never less than the weighted average of their logs:

$$\log \sum_i q_i z_i \geq \sum_i q_i \log z_i, \quad (9.15)$$

where the z_i are any set of positive numbers and the q_i are any set of non-negative weights that sum to one: $\sum_i q_i = 1$. Jensen's inequality follows from the fact that the logarithm is concave downward. For a proof see footnote 3 below.

If we make the substitution $x_i = q_i z_i$ in Eq. (9.15) we find that

$$\log \sum_i x_i \geq \sum_i q_i \log \frac{x_i}{q_i}, \quad (9.16)$$

³Suppose we have N positive numbers z_1, \dots, z_N and a further N non-negative numbers q_1, \dots, q_N that sum to one. And suppose that $f(x)$ is any linear function $f(x) = mx + c$, where m and c are constants. Then

$$f\left(\sum_i q_i z_i\right) = m \sum_i q_i z_i + c = m \sum_i q_i z_i + c \sum_i q_i = \sum_i q_i (m z_i + c) = \sum_i q_i f(z_i),$$

where we have made use of the fact that $\sum_i q_i = 1$ in the second equality. Let us choose $f(x)$ to be the linear function tangent to $\log x$ at the point $x = \sum_i q_i z_i$. Because $\log x$ is concave downward it follows that $f(x) \geq \log x$ for all positive x , with the exact equality holding at the tangent point $x = \sum_i q_i z_i$. Thus,

$$\log \sum_i q_i z_i = f\left(\sum_i q_i z_i\right) = \sum_i q_i f(z_i) \geq \sum_i q_i \log z_i,$$

and hence the result is established. Note that the same proof works for any concave-downward function—the logarithm is merely a special case.

and applying this inequality to Eq. (9.14) we get

$$\log \sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data}) \geq \sum_{\mathbf{A}} q(\mathbf{A}) \log \frac{P(\mathbf{A}, \theta | \text{data})}{q(\mathbf{A})}. \quad (9.17)$$

This result is true for any set of non-negative quantities $q(\mathbf{A})$ such that $\sum_{\mathbf{A}} q(\mathbf{A}) = 1$. It will be helpful to think of $q(\mathbf{A})$ as a (properly normalized) probability distribution over adjacency matrices \mathbf{A} .

One useful special value of $q(\mathbf{A})$ is

$$q(\mathbf{A}) = \frac{P(\mathbf{A}, \theta | \text{data})}{\sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data})}. \quad (9.18)$$

Substituting this into the right-hand side of Eq. (9.17) gives

$$\begin{aligned} \sum_{\mathbf{A}} q(\mathbf{A}) \log \frac{P(\mathbf{A}, \theta | \text{data})}{q(\mathbf{A})} &= \frac{\log \sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data})}{\sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data})} \sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data}) \\ &= \log \sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data}). \end{aligned} \quad (9.19)$$

In other words, for this special choice of $q(\mathbf{A})$, Eq. (9.17) becomes not an inequality but an exact equality. Another way to say the same thing is that over all possible choices of $q(\mathbf{A})$, the one in Eq. (9.18) is the one that maximizes the right-hand side of Eq. (9.17) (because the right-hand side is always less than or equal to the left-hand side, so its maximum possible value occurs when the two sides are equal).

Now we argue as follows. Maximizing the right-hand side of (9.17) over possible choices of $q(\mathbf{A})$ makes it equal to the left-hand side, which also makes it equal to $P(\theta | \text{data})$ by Eq. (9.13). And further maximizing $P(\theta | \text{data})$ with respect to θ will give the answer we seek—the most probable value of θ . In other words, a *double maximization* of the right-hand side of (9.17), first over $q(\mathbf{A})$ and then over θ , will give us our estimate of θ .

On the face of it, this doesn't seem like a very promising approach. We have turned what was previously a single maximization of Eq. (9.13) with respect to one quantity into a double maximization over two quantities. In fact, however, it turns out to be exactly what we need. The crucial point to notice is that when doing a double maximization it doesn't matter how we do it: we don't have to maximize over $q(\mathbf{A})$ first then over θ if we don't want to. The maximum is in the same place whatever means we use to find it. And one simple way to perform a double maximization is to maximize with respect to one thing while holding the other constant, then do the reverse, and keep on repeating the process, maximizing with respect to one then the other until we converge to the joint maximum.

This is the approach we use here. We find our solution by maximizing the right-hand side of (9.17) first with respect to $q(\mathbf{A})$ keeping θ fixed, then with respect to θ keeping $q(\mathbf{A})$ fixed, and repeating until the values stop changing. We have already seen how to maximize with respect to $q(\mathbf{A})$ —the maximum is given by Eq. (9.18). It only remains to maximize with respect to θ . Differentiating the right-hand side of (9.17) with respect to θ , while holding $q(\mathbf{A})$ fixed, and setting the result to zero gives

$$\sum_{\mathbf{A}} q(\mathbf{A}) \frac{\partial}{\partial \theta} \log P(\mathbf{A}, \theta | \text{data}) = 0. \quad (9.20)$$

(If there were more than one parameter, then there would be one equation like this for each parameter.) Now we simply solve this equation for θ . We will see an example in a moment.

The EM algorithm consists of iterating back and forth between Eqs. (9.18) and (9.20) until we reach convergence. Usually this iteration is done numerically on a computer, although there are some simple cases where it can be done by hand. In the jargon of the field, Eq. (9.18) is called the expectation step or E-step of the EM algorithm and Eq. (9.20) is called the maximization step or M-step.

We have described the EM algorithm as a method for calculating just the value of the parameter (or parameters) θ , but now we have finished deriving it we notice a beautiful feature. Inadvertently, without intending to, we have actually also calculated the ground truth \mathbf{A} . At the end of the algorithm, when we have converged to the maximum, we get a result not only for θ but also for $q(\mathbf{A})$. But Eq. (9.18) tells us that at the maximum we have

$$q(\mathbf{A}) = \frac{P(\mathbf{A}, \theta | \text{data})}{\sum_{\mathbf{A}} P(\mathbf{A}, \theta | \text{data})} = \frac{P(\mathbf{A}, \theta | \text{data})}{P(\theta | \text{data})} = P(\mathbf{A} | \text{data}, \theta). \quad (9.21)$$

In other words, the final value of $q(\mathbf{A})$ is none other than the probability of the network having adjacency matrix \mathbf{A} given the data and our value for θ . All we have to do to work out the most likely value of the ground truth is find the maximum of this quantity with respect to \mathbf{A} .

In actual fact, however, it is often better *not* to maximize with respect to \mathbf{A} but just to keep the entire probability distribution $P(\mathbf{A} | \text{data}, \theta)$. This *posterior distribution* gives us a lot of information about the ground truth: not just its most likely value but the relative probabilities of all potential ground truths. In a sense, the posterior distribution actually tells us everything we originally wanted to know. It captures the structure of the network, but also the uncertainty in that structure. If the probability distribution is strongly peaked around one value of \mathbf{A} or a small number of similar values, then we have a high degree of certainty about the structure of the network. If the distribution

is broad and spread over many different values of \mathbf{A} , then there is a lot of uncertainty.

Thus, there are two possible courses of action. The equivalent of the standard procedure for real-valued data, of reporting mean and standard deviation, would be to give the most likely value of \mathbf{A} and the parameter(s) θ that parametrize the model and hence spell out the size of the measurement error. The other possibility is to give the entire posterior distribution $P(\mathbf{A}|\theta, \text{data})$. Both are acceptable in practice. The first has the advantage of returning just a single network structure, which is often what people are expecting. The second has the advantage of giving more information about the network—sometimes much more.

9.3.5 INDEPENDENT EDGE ERRORS

Let us look at a specific example of the methods of the previous section. Suppose we make repeated measurements of the structure of a network, measuring the whole network a total of N times. Each measurement involves going through every pair of nodes to determine whether they are connected by an edge. Because of experimental error, repeated readings on the same node pair may not agree. Let us denote by E_{ij} the number of times that node pair i, j is observed to have an edge.

We also need to specify the model we will use to represent how errors are introduced into the data. There are many such models we could use, but perhaps the simplest is a model that assumes that all measurements of all node pairs are statistically independent and depend on the ground truth in the same way. Such a model can be specified using two parameters. The first, called the *true positive rate*, is the probability that we will observe an edge between two nodes where one truly exists. Let us denote this probability α . The second parameter, the *false positive rate*, is the probability that we will observe an edge where none actually exists. We will denote this β .

In terms of these parameters the probability of making N measurements on node pair i, j and observing an edge on any particular E_{ij} of those measurements is $\alpha^{E_{ij}}(1 - \alpha)^{N - E_{ij}}$ if i and j are connected by an edge in the ground-truth network, or $\beta^{E_{ij}}(1 - \beta)^{N - E_{ij}}$ if they are not. Or we can combine these two expressions and write them in the form

$$P(E_{ij}|A_{ij}, \alpha, \beta) = [\alpha^{E_{ij}}(1 - \alpha)^{N - E_{ij}}]^{A_{ij}} [\beta^{E_{ij}}(1 - \beta)^{N - E_{ij}}]^{1 - A_{ij}}. \quad (9.22)$$

Note how, by raising the first factor to the power of A_{ij} , we ensure that it appears only when there is an edge between i and j in the ground truth. Similarly, by

raising the second factor to the power of $1 - A_{ij}$ we ensure that it appears only when there is no edge.

Now the probability of the entire set of measurements, for all node pairs in the whole network, is the product of Eq. (9.22) over all pairs:

$$P(\text{data}|\mathbf{A}, \alpha, \beta) = \prod_{i < j} [\alpha^{E_{ij}}(1 - \alpha)^{N - E_{ij}}]^{A_{ij}} [\beta^{E_{ij}}(1 - \beta)^{N - E_{ij}}]^{1 - A_{ij}}. \quad (9.23)$$

This is the analog of $P(\text{data}|\mathbf{A}, \theta)$ in our earlier presentation. Note how the product is over node pairs with $i < j$, which ensures that each distinct pair gets counted only once (and pairs with $i = j$ do not get counted at all, which is correct for the simple networks we are considering here, which cannot have self-edges).

Armed with this expression, we can now apply our EM algorithm. First we employ Bayes' rule, Eq. (9.12):

$$P(\mathbf{A}, \alpha, \beta|\text{data}) = P(\text{data}|\mathbf{A}, \alpha, \beta) \frac{P(\mathbf{A})P(\alpha)P(\beta)}{P(\text{data})}. \quad (9.24)$$

As before, we will assume that the prior probabilities $P(\alpha)$ and $P(\beta)$ are uniform (in the range from zero to one, since α and β are probabilities), meaning that all values of α and β are equally likely a priori. For \mathbf{A} we could assume a uniform prior too—all networks are equally likely—but this seems like a stretch since it would mean that each node pair would have a 50% chance of being connected by an edge, which is unrealistic for most networks. As we have seen, most real-world networks are very sparse, with the average probability of an edge being much less than 50%. So instead let us adopt a prior in which each edge appears with some other probability ρ , so that

$$P(\mathbf{A}|\rho) = \prod_{i < j} \rho^{A_{ij}}(1 - \rho)^{1 - A_{ij}}. \quad (9.25)$$

This introduces a third parameter into our model, so Eq. (9.24) becomes

$$P(\mathbf{A}, \alpha, \beta, \rho|\text{data}) = P(\text{data}|\mathbf{A}, \alpha, \beta) \frac{P(\mathbf{A}|\rho)P(\rho)P(\alpha)P(\beta)}{P(\text{data})}, \quad (9.26)$$

where we will also assume that the prior $P(\rho)$ is uniform.

Putting together Eqs. (9.23), (9.25), and (9.26), we now have

$$P(\mathbf{A}, \alpha, \beta, \rho|\text{data}) = \frac{1}{P(\text{data})} \prod_{i < j} [\rho \alpha^{E_{ij}}(1 - \alpha)^{N - E_{ij}}]^{A_{ij}} [(1 - \rho)\beta^{E_{ij}}(1 - \beta)^{N - E_{ij}}]^{1 - A_{ij}}. \quad (9.27)$$

We will use our EM algorithm to maximize this expression and calculate the most likely values of the parameters with the data fixed at their observed values. Note that since the data are fixed, the prior probability $P(\text{data})$ is also fixed, so it has no effect on the position of the maximum.

Substituting (9.27) into the ‘‘E-step’’ equation of the EM algorithm, Eq. (9.18), we get

$$\begin{aligned}
 q(\mathbf{A}) &= \frac{\prod_{i<j} [\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}}]^{A_{ij}} [(1-\rho)\beta^{E_{ij}}(1-\beta)^{N-E_{ij}}]^{1-A_{ij}}}{\sum_{\mathbf{A}} \prod_{i<j} [\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}}]^{A_{ij}} [(1-\rho)\beta^{E_{ij}}(1-\beta)^{N-E_{ij}}]^{1-A_{ij}}} \\
 &= \prod_{i<j} \frac{[\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}}]^{A_{ij}} [(1-\rho)\beta^{E_{ij}}(1-\beta)^{N-E_{ij}}]^{1-A_{ij}}}{\sum_{A_{ij}=0,1} [\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}}]^{A_{ij}} [(1-\rho)\beta^{E_{ij}}(1-\beta)^{N-E_{ij}}]^{1-A_{ij}}} \\
 &= \prod_{i<j} Q_{ij}^{A_{ij}} (1-Q_{ij})^{1-A_{ij}}, \tag{9.28}
 \end{aligned}$$

where

$$Q_{ij} = \frac{\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}}}{\rho\alpha^{E_{ij}}(1-\alpha)^{N-E_{ij}} + (1-\rho)\beta^{E_{ij}}(1-\beta)^{N-E_{ij}}}. \tag{9.29}$$

In other words, the posterior distribution over networks \mathbf{A} factors into a product of terms, one for each node pair, with each one depending on just a single quantity Q_{ij} . The value of Q_{ij} represents the probability, given the observed measurements, that there is an edge between nodes i and j . One can think of the Q_{ij} as the elements of an $n \times n$ symmetric matrix \mathbf{Q} that is a natural generalization of the adjacency matrix. If Q_{ij} is either zero or one then it behaves like an ordinary adjacency matrix element—there either is or is not an edge between nodes i and j . For values in between it describes both the observations and their measurement error. For instance, if $Q_{ij} = 0.9$ then there is a strong chance that i and j are connected by an edge in the network, but there is still a 10% chance that they are not, which represents our experimental uncertainty. Conversely if $Q_{ij} = 0.1$ then probably there is no edge between i and j , but again there is a 10% chance that we are wrong. In this way both the measured result and the error are captured in a single number, by contrast with the traditional approach for real-valued data, which requires two numbers, a mean and a standard deviation, to represent the result and the error.⁴

The other half of our EM algorithm is the M-step of Eq. (9.20). Because there are three parameters in our model, α , β , and ρ , Eq. (9.20) becomes three

⁴Formally, this is because the edges have a Bernoulli distribution, not a normal distribution, and the Bernoulli distribution is completely specified by one (probability) parameter, by contrast with the normal distribution, which requires two parameters, a mean and a standard deviation.

equations here, one for differentiation with respect to each of the parameters. Using Eq. (9.27), performing the derivatives, and employing (9.28) for $q(\mathbf{A})$, we find that

$$\alpha = \frac{\sum_{i<j} E_{ij} Q_{ij}}{N \sum_{i<j} Q_{ij}}, \quad \beta = \frac{\sum_{i<j} E_{ij} (1 - Q_{ij})}{N \sum_{i<j} (1 - Q_{ij})}, \quad \rho = \frac{1}{\binom{n}{2}} \sum_{i<j} Q_{ij}. \quad (9.30)$$

Iterating Eqs. (9.29) and (9.30) repeatedly now gets us the full solution for our parameters α , β , and ρ , and the edge probabilities Q_{ij} . Typically, we would perform the calculation numerically, starting from any reasonable set of initial values and iterating until the results converge. It's usually easier to choose initial values of the parameters and then apply Eq. (9.29) to calculate Q_{ij} , rather than choosing initial values for the Q_{ij} , simply because there are so many of the latter, whereas there are only three parameters.

9.3.6 EXAMPLE

As an example of the methods of the previous section we consider a data set generated by Eagle and Pentland in what they call a “reality mining” experiment [154, 155]. This experiment, carried out in 2004 and 2005, was of a type discussed previously in Section 4.4, in which the experimenters measured social networks of face-to-face interaction using mobile phones. A group of 96 university students were given phones equipped with special software that recorded when two phones were in close physical proximity (a few meters or less) using Bluetooth radio technology. The catch with this kind of experiment is that observed proximity doesn't necessarily guarantee that people actually have a social connection. They might just pass each other on the street or eat at the same restaurant. On the other hand, such chance connections will typically be only occasional and sporadic, whereas people who have a real social connection will probably be in proximity on a more regular basis. Thus proximity, as recorded in this experiment, does reflect social connections, but it is an error-prone measurement of the social network. This is exactly the kind of data that the method above is designed to deal with.

Here we'll take a subset of the data collected in the experiments, representing proximity measurements from eight consecutive Wednesdays in March and April of 2005. For each of the eight days we record for every pair of participants whether they were, or were not, observed in proximity at any time during that day. Thus in this case the number of measurements of the network is $N = 8$ and the values of E_{ij} can run from 0 to 8 (and all of those values do in fact occur in the data). The reason for choosing to look at the same single day each week is that there is a considerable amount of weekly variation in the data. There

Be careful to distinguish here between N , which is the number of times we measure the network, and n , which is the number of nodes in the network.

In practice one often uses random initial values, although other choices may be appropriate, particularly if one has some inkling beforehand of what the correct answers should be.

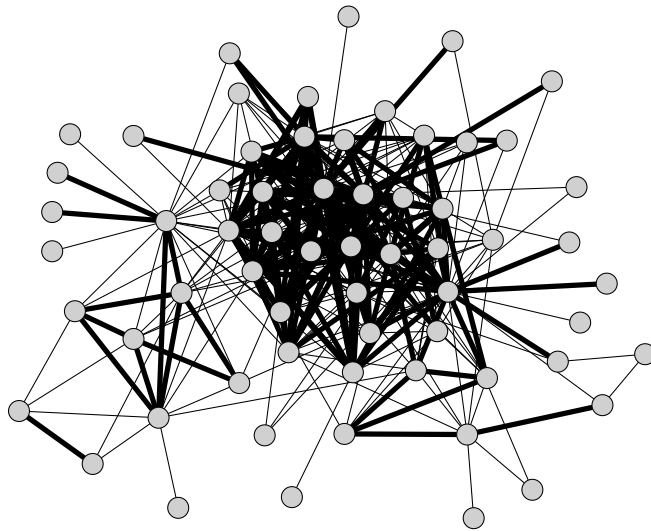


Figure 9.2: A social network deduced from proximity data. The thicknesses of the lines indicate our certainty Q_{ij} about whether an edge exists or not. Thicker lines indicate a higher probability that there really is an edge between the nodes in question.

are, for example, many fewer interactions between participants at the weekend than on weekdays. By looking at the same day each week we filter out this weekly variation.

Taking the measurements, feeding them into Eqs. (9.29) and (9.30), and iterating the equations to convergence, gives us values $\alpha = 0.4242$, $\beta = 0.0043$, and $\rho = 0.0335$ for the three model parameters. The small value of ρ tells us that the network is quite sparse, but otherwise it is not of particular interest. The small value of β tells us that the false positive rate is low for this experiment, less than 1%. On the other hand, the value $\alpha = 0.42$ means that the false negative rate (which is $1 - \alpha$) is relatively high, over 50%. This is not necessarily a bad thing. It simply means that not everyone has contact with all of their acquaintances every day, which seems plausible.

Our calculation also gives us the values of the Q_{ij} , which are shown in Fig. 9.2. The thickness of the edges in the figure represent the probabilities of the edges and, as we can see, there is a core of about twenty nodes in the network that are connected with high probability, plus a number of peripheral nodes with weaker connections.

Figure 9.3 shows a plot of the edge probability, Eq. (9.29), as a function of the number of proximity observations E_{ij} (which, as we have said, ranges from 0 to 8). As the figure shows, if $E_{ij} = 0$ or 1 the value of Q_{ij} is small—less than 10%. But if E_{ij} is 2 or more then Q_{ij} jumps to over 90%. In other words, according to our best estimate, if two people are observed in proximity only once during the experiment it is probably a false alarm. They were just passing strangers. But if they are seen together twice or more then they probably know each other.

9.3.7 ESTIMATION OF OTHER QUANTITIES

Once one has estimated the ground-truth network \mathbf{A} , one can also calculate other quantities of interest, such as degrees, centrality measures, path lengths, clustering coefficients, and so forth. Better still, armed with the posterior distribution over networks $P(\mathbf{A}|\text{data}, \theta)$, one can (at least in principle) calculate the entire distribution over any other quantity that depends on \mathbf{A} . Suppose that $X(\mathbf{A})$ is some quantity of interest that depends on the structure of the network. Then the probability that it takes a particular value x is given by

$$P(X = x) = \sum_{\mathbf{A}} \delta_{x, X(\mathbf{A})} P(\mathbf{A}|\text{data}, \theta), \quad (9.31)$$

where δ_{ij} is the Kronecker delta. Alternatively, we can calculate the mean or expected value of X as an average over \mathbf{A} :

$$\langle X \rangle = \sum_{\mathbf{A}} X(\mathbf{A}) P(\mathbf{A}|\text{data}, \theta), \quad (9.32)$$

or the standard deviation

$$\sigma_X^2 = \sum_{\mathbf{A}} [X(\mathbf{A}) - \langle X \rangle]^2 P(\mathbf{A}|\text{data}, \theta). \quad (9.33)$$

For example, suppose we are interested in the degree of node i , which we can write as $k_i = \sum_j A_{ij}$ (Eq. (6.12)). Then the expected value of this degree is

$$\begin{aligned} \langle k_i \rangle &= \sum_{\mathbf{A}} \sum_j A_{ij} P(\mathbf{A}|\text{data}, \theta) = \sum_j \sum_{\mathbf{A}} A_{ij} P(\mathbf{A}|\text{data}, \theta) \\ &= \sum_j \sum_{A_{ij}=0,1} A_{ij} P(A_{ij}|\text{data}, \theta) = \sum_j P(A_{ij} = 1|\text{data}, \theta). \end{aligned} \quad (9.34)$$

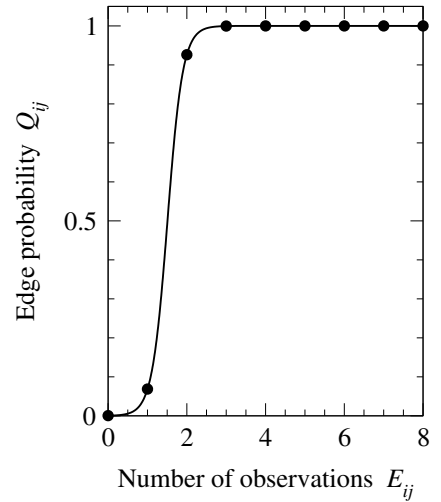


Figure 9.3: Edge probabilities in the “reality mining” experiment. The estimated probability Q_{ij} of an edge between a pair of nodes i, j as a function of the number of times such an edge was observed in the data.

The quantity $P(A_{ij} = 1 | \text{data}, \theta)$ is simply the probability that there is an edge between nodes i and j , which is the same as the quantity we called Q_{ij} in the independent edge model of Section 9.3.5. So in that case we have $\langle k_i \rangle = \sum_j Q_{ij}$, a simple generalization of the standard expression $k_i = \sum_j A_{ij}$ for node degree in a network without experimental uncertainty.

In many cases, calculating an expected value or distribution for a quantity of interest, rather than just giving the value for the most likely network structure, is a better way to go, since it captures not only the most likely structure but also other competing structures that might have a probability almost as high.

9.3.8 OTHER ERROR MODELS

The independent edge error model of Section 9.3.5 is simple and intuitive. It is a good starting point for error analysis in many networks. Unlike the situation for real-valued data, however, where one model, the Gaussian model, is used in almost all cases, network data call for different models in different circumstances [361]. One can imagine a number of variants on the independent edge model. One could have different true- and false-positive rates in different parts of the network, for instance, or rates that depend on other properties of nodes or edges. One could allow for correlations between edges or different numbers of measurements of different node pairs. One could also consider weighted networks in which edges can have different strengths. In a social network like the reality mining example of the previous section, for instance, one could assume that rather than just two possible states for node pairs, not acquainted and acquainted, one has instead three—not acquainted, somewhat acquainted, and strongly acquainted. (Four or more levels would be straightforward to define too.)

There are also cases where a completely different model is needed. For instance, we mentioned in Section 9.2 the issues with fixed choice social network surveys, in which there is a maximum number of contacts a person can report. These features could be incorporated into an error model too. In principle all one need do is write down the probability $P(\text{data} | \mathbf{A}, \theta)$ that a particular set of observations will occur given the underlying ground truth, although in practice doing so may not be easy and the probabilities may have to be calculated or approximated using numerical methods. A selection of more complicated error models for social network data have been discussed by Butts [90], who focused particularly on the issue of participant accuracy in social network surveys. In these models it is assumed that the primary source of error in the network is variation in the reliability of responses given by participants, some reporting their social contacts accurately and some not. This introduces a large set of

additional parameters quantifying individual reliability and the models can become quite complex.

A particular problem with the independent edge model of the previous section is that it requires us to measure the structure of the network more than once in order to apply Eqs. (9.29) and (9.30). In a sense this is inevitable, since it is well known that you cannot make an error estimate from a single data point, but it is also problematic since in practice it is quite rare for experimenters to make multiple measurements of the same network. One case in which we do have multiple measurements is in studies of friendship networks or similar social networks using interviews or questionnaires. In such studies one effectively has two observations of each friendship, from the point of view of each of the two people involved, and one can make an estimate of error rates from how often people agree about whether they are friends [90,362].

There are also certain circumstances in which it is possible to estimate errors from just a single network measurement. In particular, if the edges of a network are *correlated*, then we may be able to estimate errors. To give an (admittedly unrealistic) example, suppose there are two edges in a network that are perfectly correlated, meaning that we know them to be either both present in the network or both absent (but we don't know which). In that case, if we make a single measurement of the structure of the whole network, including these two edges, then we are in effect measuring the same quantity twice—the two edges represent the same measured quantity because they always have the same value. Thus we have two data points for the given edge and hence we can make an error estimate. In real life edges are rarely perfectly correlated in this way, but even partial correlations can be used to make error estimates. A number of versions of this idea have been pursued in the literature—see Refs. [109] and [225], for example—although the theory has not yet been fully explored and there is room for further work.

9.4 CORRECTING ERRORS

We have seen one way to cope with errors in network data, by estimating their size and the effect they will have on other measurements we make. Another approach is to try and fix the errors themselves, to improve the quality of the data. Network data are unusual in that we often have a lot of insight into the likely structure of the networks we are examining, which in some cases allows us to guess with some accuracy where the errors are and then set them right. There are two situations in particular in which this approach is commonly used: link prediction and node disambiguation.

9.4.1 LINK PREDICTION

A good example of error correction in networks is *link prediction*, meaning the identification of false negative edges—those that are erroneously missing from a network. One could in principle also try to identify false positive edges, those that are erroneously included in the network, although this is done less often.

Suppose we have an observed network and we believe that there may be some edges missing from it. That is, there are edges in the real, ground-truth network that do not appear in the empirical data. Given the data we have, can we make a guess about what is missing? This is a rather different problem from the error estimation questions of previous sections. Usually when doing link prediction we have only a single measurement of the entire network and our predictions are based on assumptions about correlations between edges. In other words, by contrast for instance with the approach of Section 9.3.5, we explicitly assume that the edges are not statistically independent and use this to make guesses about where the data are in error.

A number of simple link prediction methods have been shown to work well, at least in some circumstances. These are really no more than rules-of-thumb and are not based on particular models or derivations. Rather, one simply guesses a strategy that might work and then tests it out to see how well it fares in practice. The standard way to perform the testing is *cross-validation*, in which one takes a known network, removes some edges from it, then tests to see whether the proposed prediction method is able to identify the missing edges. In practice, no method is able to predict all of the missing edges all of the time. Rather, most methods return a list of node pairs ordered from most to least likely to be connected by a missing edge. Even the node pairs at the top of the list, however, may not have a very high likelihood of being correct. In this game, merely doing better than chance is considered a win. If one were to guess at random where a missing edge was, one would have a probability of about $1/\binom{n}{2}$ of guessing right, since this is the number of pairs to choose between.⁵ For large n this is a very small number, so a link prediction method does not have to succeed with very high probability to be better than a random guess.

Examples of the kinds of strategies used include:

1. *Shortest path distance*: Node pairs are considered more likely to be connected by a missing edge if the graph distance between them is small.
2. *Number of common neighbors*: Node pairs with a larger number of common

⁵To be exact the probability is $1/[\binom{n}{2} - m]$ where m is the number of observed edges. But in the common case of a sparse network m is much less than $\binom{n}{2}$ and hence can be ignored.

neighbors are considered more likely to be connected than those with fewer.

3. *Node degrees*: Nodes with high degree are considered more likely to be connected. Typically, one looks at the product $k_i k_j$ of degrees of node pairs, which is proportional to the probability of an edge in the so-called configuration model—see Eq. (12.2) and the preceding discussion.
4. *Node similarity*: Similar nodes are considered more likely to be connected than less similar ones. Similarity can be measured using any of the measures described in Section 7.6, such as cosine similarity or the Jaccard coefficient.

For example, Liben-Nowell and Kleinberg [301] tested nine previously proposed methods of link prediction on a set of scientific collaboration networks, in effect trying to predict which scientists were likely to have collaborated even if those collaborations were unobserved. Performing a cross-validation test of the kind described above, they measured the factor by which each method improved over a random guess. They found improvements of up to a factor of about 50 for some methods, which sounds impressive, although one must remember that the probability of a random guess being right is small in the first place—typically around 0.2% in this study—so the test is setting a fairly low bar for success. Even with a factor of 50 improvement, the overall probability of making a correct prediction is still only about $50 \times 0.2\% = 10\%$, meaning that 90% of predictions are still wrong.

Perhaps the best way to think about calculations like this is as a guide to further experiments: they cannot tell us for sure where the missing edges are, but they can give us a hint where to look. With a probability of 0.2% of finding an edge on a random guess we are going to have to look through about 500 pairs of nodes before we have even one success. With a probability of 10% on the other hand, we are only going to have to look through about ten pairs. That could make a big difference, particularly in cases, such as biological networks, where the amount of work needed to establish the existence of even a single edge can be substantial.

There are also more rigorous ways of predicting missing links, although they tend to be more complicated to implement than the simple heuristics described above. One approach is to fit the observed network to a network model and then use the model to calculate the probabilities of appearance of individual edges. A typical example of this approach involves using the “stochastic block model,” a model that we study in detail in Sections 12.11.6 and 14.4.1. For instance, Guimerà and Sales-Pardo [225] used the stochastic block model combined with a Bayesian fitting technique based on Monte Carlo sampling to perform link prediction on a number of networks, including a

protein–protein interaction network and a network of airline routes. Methods like this can outperform simpler heuristics by a significant margin, particularly on certain types of networks [109]. Whether this improvement justifies their substantial additional complexity is a matter of individual opinion.

9.4.2 NODE DISAMBIGUATION

Node disambiguation is the (rather ungainly) name given to the process of trying to discern when two nodes in a network are actually duplicates of one another, or conversely when two nodes have been inadvertently combined and should be separated again. As discussed in Section 9.1, these issues arise particularly in social networks where individuals are identified by name alone. In collaboration and coauthorship networks, for instance, if authors are identified only by name then we may well inadvertently confuse two people with the same name, or mistake one person for two if they give their name differently on different occasions.

Node disambiguation, also sometimes called *entity resolution*, can be done in various ways [127, 173, 430, 444], but one of the nice features of the problem is that we can often make use of the network itself to give us pointers to how to proceed. In a collaboration network, for instance, two nodes with similar names (“J. Doe” and “Jane Doe” say) might or might not represent the same person. But if we observe that they have a lot of the same collaborators then it makes it much more likely that they are in fact one person. In addition to this kind of network-driven inference one may also be able to make use of additional data associated with nodes, such as geographical location. If we observe that our J. Doe and Jane Doe both work at the same institution, for instance, that too increases the chances that they are the same person.

Typically, node disambiguation methods start by assuming that every instance of a node in the observed network is distinct. Every author of every paper in a scientific coauthorship network, for instance, is assumed, initially, to be a new person we have not seen before. The goal of disambiguation is then to amalgamate nodes that represent the same person, in order to create the final network.

For instance, Ferreira *et al.* [173] review a range of different methods for disambiguation in coauthorship networks based on network measures of author similarity. One takes pairs of authors with similar names, calculates a measure such as cosine similarity (Section 7.6.1), then amalgamates authors if their similarity rises above a certain threshold. In effect, this amalgamates authors if they have many of the same collaborators, which, as described above, seems like a sensible approach. Other measures of similarity have been tried as well,

such as the Jaccard coefficient (Section 7.6.1) and more complicated measures based on machine learning techniques.

EXERCISES

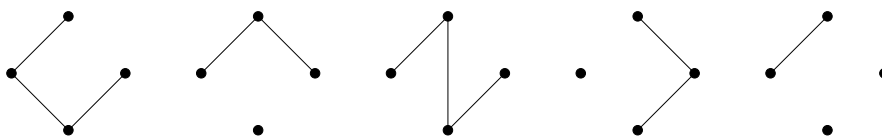
9.1 Suggest some potential sources of error in measurements of the structure of the following networks:

- A scientific coauthorship network assembled from a database of papers.
- A web network of web pages at a single university, assembled by using an automated web crawler.
- A metabolic network.
- A social network of who is friends with whom at a large company, assembled using questionnaires.
- A network representation of an electrical power grid.

9.2 Suppose we draw n independent random reals x in the range $0 \leq x < \infty$ from the (properly normalized) exponential probability density $P(x) = \mu e^{-\mu x}$.

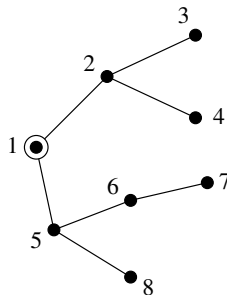
- Write down the likelihood (i.e., the probability density) that we draw a particular set of n values x_i (where $i = 1 \dots n$) for a given value of the exponential parameter μ .
- Hence find a formula for the best (meaning the maximum-likelihood) estimate of μ given a set of observed values x_i .

9.3 A small (4-node) network is measured five times, but the measurements are unreliable so the observed structure is different each time. Here are the five structures observed for the network:



Applying the method of Section 9.3.5, calculate the probability Q_{ij} that each of the six possible edges exists. (Hint: You will probably need to write a computer program to solve this problem.)

9.4 Recall that the structure of the Internet is measured by reconstructing paths from traceroute or BGP data (see Sections 2.1.1 and 2.1.2). The real Internet is a large network but for illustrative purposes consider this small example of eight nodes:



Here we have measured the shortest path from node 1 on the left to all others (or one such path to each node if there are several). The edges shown are the ones that make up these paths, and we believe them to be reliable measurements, but they may not constitute the whole network—there may be some additional edges that are not observed and don't appear in the figure.

- a) Given that we know that the observed edges constitute shortest paths from node 1 to all other nodes, about which of the remaining pairs of nodes in the network are we now uncertain? That is, for which pairs of nodes can we not tell whether they are connected by an edge?
- b) Give a general rule for determining, for any given node pair in this network, whether they are (i) definitely connected, (ii) definitely not connected, or (iii) their connection status is unknown.

One could use such a rule to formulate an error model of the kind discussed in Section 9.3.3 for Internet data.

9.5 The false positive rate β defined in Section 9.3.5 is the probability of erroneously observing an edge where none exists. Arguably a more useful measure, however, is the *false discovery rate*, which is the probability that an actual observed edge is itself a false positive, which is not the same thing.

Suppose we measure a network once, producing an observed adjacency matrix \mathbf{O} with elements O_{ij} . Then the standard false positive rate within the independent edge model of Section 9.3.5 is $\beta = P(O_{ij} = 1 | A_{ij} = 0, \alpha, \beta, \rho)$. (It does not matter which pair of nodes i, j we look at. By hypothesis all pairs have the same false positive rate.)

The probability that an observed edge is a false positive, on the other hand, is given by $P(A_{ij} = 0 | O_{ij} = 1, \alpha, \beta, \rho)$. Using Bayes' rule, we can write this probability as

$$P(A_{ij} = 0 | O_{ij} = 1, \alpha, \beta, \rho) = P(O_{ij} = 1 | A_{ij} = 0, \alpha, \beta, \rho) \frac{P(A_{ij} = 0 | \alpha, \beta, \rho)}{P(O_{ij} = 1 | \alpha, \beta, \rho)}.$$

a) Show that

$$P(O_{ij} = 1|\alpha, \beta, \rho) = P(O_{ij} = 1|A_{ij} = 1, \alpha, \beta, \rho)P(A_{ij} = 1|\alpha, \beta, \rho) \\ + P(O_{ij} = 1|A_{ij} = 0, \alpha, \beta, \rho)P(A_{ij} = 0|\alpha, \beta, \rho).$$

b) Hence show that the probability of an observed edge being a false positive is

$$P(A_{ij} = 0|O_{ij} = 1, \alpha, \beta, \rho) = \frac{\beta(1 - \rho)}{\alpha\rho + \beta(1 - \rho)}.$$

c) Calculate the value of this probability for the “reality mining” example of Section 9.3.6 using the values of α , β , and ρ given on page 294.

d) You should find that the probability you calculate is considerably larger than the false negative rate β . Explain briefly why this is.

9.6 For some networks, observations of edges are reliable but observations of non-edges are not. Academic coauthorship networks (Section 4.5) provide an example. If we observe a paper written by two particular individuals then it is a safe bet that they did actually coauthor a paper. But if we do not observe such a paper, then it does not guarantee that no such paper exists. We might just not have found it, or it might be written but not published yet.

Modify the error model used in Section 9.3.5 to this case of reliable edges but unreliable non-edges and derive the equivalent of Eqs. (9.29) and (9.30).

CHAPTER 10

THE STRUCTURE OF REAL-WORLD NETWORKS

A discussion of some of the features and patterns that are revealed when we apply the concepts developed in previous chapters to observed networks

IN PREVIOUS chapters of this book we have looked at the various types of natural and man-made networks and methods for determining their structure (Chapters 2 to 5), the mathematical techniques used to represent and quantify networks (Chapters 6 and 7), and the computer algorithms and statistical methods necessary for practical analysis of today's large network data sets (Chapters 8 and 9). In this chapter we combine what we have learned so far, applying our theoretical ideas, measures, and methods to empirical network data to determine what networks look like in the real world.

As we will see, there are a number of distinctive patterns in the structure of real-world networks that recur over and over again and can have a profound effect on the way networked systems behave. Among other things, we discuss component sizes, path lengths and the small-world effect, degree distributions and power laws, clustering coefficients, and network correlations and assortative mixing. These are some of the key concepts of the field.

10.1 COMPONENTS

We begin our discussion of the structure of real-world networks with a look at their component structure. (See Section 6.12 for definitions and discussion of

	Network	Type	n	m	c	S	ℓ	α	C	C_{WS}	r	Ref(s).
Social	Film actors	Undirected	449 913	25 516 482	113.43	0.980	3.48	2.3	0.20	0.78	0.208	20,466
	Company directors	Undirected	7 673	55 392	14.44	0.876	4.60	–	0.59	0.88	0.276	131,369
	Math coauthorship	Undirected	253 339	496 489	3.92	0.822	7.57	–	0.15	0.34	0.120	133,219
	Physics coauthorship	Undirected	52 909	245 300	9.27	0.838	6.19	–	0.45	0.56	0.363	347,349
	Biology coauthorship	Undirected	1 520 251	11 803 064	15.53	0.918	4.92	–	0.088	0.60	0.127	347,349
	Telephone call graph	Undirected	47 000 000	80 000 000	3.16			2.1				10,11
	Email messages	Directed	59 812	86 300	1.44	0.952	4.95	1.5/2.0		0.16		156
	Email address books	Directed	16 881	57 029	3.38	0.590	5.22	–	0.17	0.13	0.092	364
	Student dating	Undirected	573	477	1.66	0.503	16.01	–	0.005	0.001	–0.029	52
Sexual contacts	Undirected	2 810					3.2				304,305	
Information	WWW nd. edu	Directed	269 504	1 497 135	5.55	1.000	11.27	2.1/2.4	0.11	0.29	–0.067	16,41
	WWW AltaVista	Directed	203 549 046	1 466 000 000	7.20	0.914	16.18	2.1/2.7				84
	Citation network	Directed	783 339	6 716 198	8.57			3.0/–				404
	Roget's Thesaurus	Directed	1 022	5 103	4.99	0.977	4.87	–	0.13	0.15	0.157	272
	Word co-occurrence	Undirected	460 902	16 100 000	66.96	1.000		2.7		0.44		146,175
Technological	Internet	Undirected	10 697	31 992	5.98	1.000	3.31	2.5	0.035	0.39	–0.189	102,168
	Power grid	Undirected	4 941	6 594	2.67	1.000	18.99	–	0.10	0.080	0.003	466
	Train routes	Undirected	587	19 603	66.79	1.000	2.16	–		0.69	–0.033	425
	Software packages	Directed	1 439	1 723	1.20	0.998	2.42	1.6/1.4	0.070	0.082	–0.016	352
	Software classes	Directed	1 376	2 213	1.61	1.000	5.40	–	0.033	0.012	–0.119	453
	Electronic circuits	Undirected	24 097	53 248	4.34	1.000	11.05	3.0	0.010	0.030	–0.154	174
	Peer-to-peer network	Undirected	880	1 296	1.47	0.805	4.28	2.1	0.012	0.011	–0.366	6,409
Biological	Metabolic network	Undirected	765	3 686	9.64	0.996	2.56	2.2	0.090	0.67	–0.240	252
	Protein interactions	Undirected	2 115	2 240	2.12	0.689	6.80	2.4	0.072	0.071	–0.156	250
	Marine food web	Directed	134	598	4.46	1.000	2.05	–	0.16	0.23	–0.263	245
	Freshwater food web	Directed	92	997	10.84	1.000	1.90	–	0.20	0.087	–0.326	321
	Neural network	Directed	307	2 359	7.68	0.967	3.97	–	0.18	0.28	–0.226	466,470

Table 10.1: Basic statistics for a number of networks. The properties measured are: type of network, directed or undirected; total number of nodes n ; total number of edges m ; mean degree c ; fraction of nodes in the largest component S (or the largest weakly connected component in the case of a directed network); mean distance between connected node pairs ℓ ; exponent α of the degree distribution if the distribution follows a power law (or “–” if not; in/out-degree exponents are given for directed networks); clustering coefficient C from Eq. (7.28); clustering coefficient C_{WS} from the alternative definition of Eq. (7.31); and the degree correlation coefficient r from Eq. (7.64). The last column gives the citation(s) for each network in the References. Blank entries indicate unavailable data.

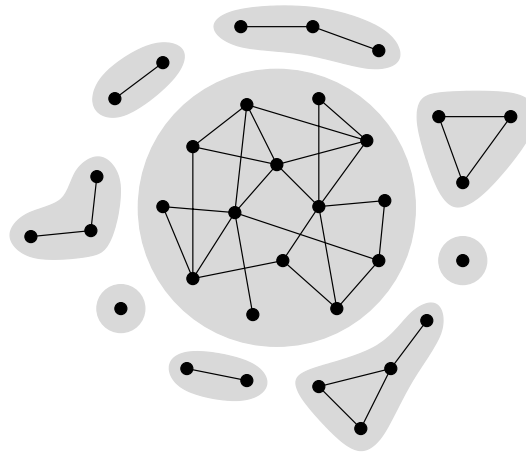


Figure 10.1: Components in an undirected network. In most undirected networks there is a single large component occupying a significant fraction of the network, along with a number of small components, typically consisting of only a handful of nodes each.

components in networks.)

In an undirected network, we typically find that there is a large component that fills most of the network—usually more than half and not infrequently over 90% of the nodes—while the rest of the network is divided into a large number of small components disconnected from the rest. The situation is sketched in Fig. 10.1. (The large component is often referred to as the “giant component,” although this is a slightly sloppy usage. As discussed in Section 11.5, the words “giant component” have a specific meaning in network theory and are not precisely synonymous with “largest component.” In this book we will be careful to distinguish between “largest” and “giant.”)

A typical example of this kind of behavior is the network of film actors discussed in Section 4.5. In this network the nodes represent actors in movies and there is an edge between two actors if they have ever appeared in the same movie. In a version of the network from May 2000 [369], it was found that 440 971 out of 449 913 actors were connected together in the largest component, or about 98%. Thus just 2% of actors were not part of the largest component.

Table 10.1 summarizes the properties of many of the networks discussed in this chapter, and gives, among other things, the size S of the largest component in each case as a fraction of total network size. (For the directed networks in the table it is the size of the largest weakly connected component that is quoted,

i.e., the largest component size when one simply ignores the directions of the edges—see Section 6.12.1. Component sizes in directed networks are discussed further in the following section.) As we can see from the table our figure for the actor network is quite typical for the networks listed and not unusually large.

As the table also shows, there are quite a few networks for which the largest component fills the entire network so that $S = 1$, meaning that the network has only a single component and no smaller components. In the cases where this happens there is usually a good reason. Take the Internet, for example. The Internet is a communication network—its reason for existence is to provide connections between its nodes. There would be little point in being a part of the Internet if you are not part of its largest component, since that would mean that you are disconnected from and unable to communicate with almost everyone else. Thus there is a strong pressure on every node that is connected to the Internet at all to belong to its largest component and thus for the largest component to fill the entire network.

In some other cases the largest component fills the network because of the way the network is measured. The first web network listed in Table 10.1 for instance, the network for the `nd.edu` web domain, is derived from a single web crawl (see Section 3.1). Since a crawler can only find a web page if that page is linked to by another page, it follows automatically that all pages found by a single crawl will be connected in a single component. The full network of the entire World Wide Web certainly has many components, but the subset captured by a single crawl has only one. The second web network listed in the table, the “AltaVista” network, was assembled using several web crawls starting from different locations, and this network does have more than one component, so $S < 1$.

Can a network have two or more large components that fill a sizable fraction of the entire network? Usually the answer to this question is no. We will study this point in more detail in Section 11.5.1, but the basic argument is this. If we had a network of n nodes that was divided into two large components of about $\frac{1}{2}n$ nodes each, then there would be $\frac{1}{4}n^2$ possible pairs of nodes such that one node was in one large component and the other node in the other large component. If there is an edge between *any* of these pairs of nodes, then the two components are joined together and are in fact just one component. For example, in our network of movie actors, with half a million nodes, there would be about 50 billion node pairs with one node in each of the two halves, and only one of those pairs would have to be connected by an edge to join the two large components into one. Except in very special cases, it is highly unlikely that not a single such pair would be connected, and hence also highly unlikely that we will have two large components.

And what about networks with no large component? It is certainly possible for networks to consist only of small components, small groups of nodes connected among themselves but not connected to the rest of the world. An example would be the network of immediate family ties, in which two people are considered connected if they are family members living under the same roof. Such a network is clearly formed of many small components consisting of individual families, with no large component at all. In practice, however, situations like this arise rather infrequently in the study of networks for the simple reason that people don't usually bother to represent such situations by networks at all. Network representations of systems are normally only useful if most of the network is connected together. If a network is so sparse as to be made only of small components, then there is little to be gained by applying techniques like those described in this book. Thus, essentially all of the networks we will be looking at do contain a large component (and certainly all those in Table 10.1, although for some of them the size of that component has not been measured and the relevant entry in the table is blank).

So the picture we have of the component structure of most networks is that of Fig. 10.1, of a large component filling most of the network, sometimes all of it, and perhaps some other small components that are not connected to the bulk of the network.

10.1.1 COMPONENTS IN DIRECTED NETWORKS

As discussed in Section 6.12, the component structure for directed networks is more complicated than for undirected ones. Directed networks have weakly and strongly connected components. The weakly connected components correspond closely to the concept of a component in an undirected network, and the typical situation for weakly connected components is similar to that for undirected networks: there is usually one large weakly connected component plus, optionally, some other small ones. Figures for the fractional sizes of the largest weakly connected components in several directed networks are given in Table 10.1.

A strongly connected component, as described in Section 6.12, is a set of nodes such that each can reach and is reachable from all others in the set along a directed path. As with weakly connected components, there is typically one large strongly connected component in a directed network and a selection of small ones. The largest strongly connected component of the World Wide Web, for instance, fills about a quarter of the network [84].

Associated with each strongly connected component is an out-component (the set of all nodes that can be reached along a directed path from any starting

point in the strongly connected component) and an in-component (the set of nodes from which the strongly connected component can be reached). By their definition, in- and out-components are supersets of the strongly connected component to which they belong and if there is a large strongly connected component in a network then the portions of the corresponding in- and out-components that lie outside the strongly connected component will often also be large. In the Web, for example, the portions of the in- and out-components that lie outside the largest strongly connected component each also occupy about a quarter of the network [84].

Each of the small strongly connected components in a directed network will have its own in- and out-components also. Often these will themselves be small, but they need not be. It can happen, for instance, that a small strongly connected component \mathcal{C} is connected by a directed path to the large strongly connected component, in which case \mathcal{C} has a large out-component despite being small itself. Note that the large strongly connected component can be reachable from many small components in this way—the out-components of different strongly connected components can overlap in directed networks and there are usually many nodes that belong to more than one out-component. Similar arguments apply for in-components as well.

The overall picture for a directed network can be represented using the “bow tie” diagram introduced by Broder and co-workers [84]. In Fig. 10.2 we show the bow tie for the case of the World Wide Web, including percentages (from Ref. [84]) for the fraction of the network occupied by its different parts.¹

Not all directed networks have a large strongly connected component. In particular, an acyclic network has no strongly connected components of size greater than one since if two nodes belong to the same strongly connected component then by definition there exists a directed path in both directions between them, and hence there is a cycle from one node to the other and back. Real-life networks are not usually perfectly acyclic, but some, such as citation networks, are approximately so (see Section 3.2). Such networks typically have a few small strongly connected components of two or perhaps three nodes each, but no large ones.

¹The study of Ref. [84], which was published in 2000, is now quite old and the web network has grown and changed significantly since it was performed. It is quite possible that the relative sizes of the various parts of the network have changed too, perhaps considerably, in the years since this paper was published.

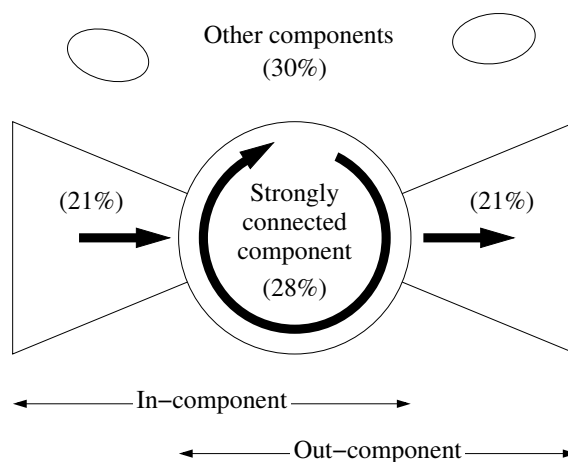


Figure 10.2: The “bow tie” diagram of components in a directed network. The typical directed network consists of one large strongly connected component and many small ones, each with an in-component and an out-component. Note that by definition each in-component includes the corresponding strongly connected component as a subset, as does each out-component. The largest strongly connected component and its in- and out-components typically occupy a significant fraction of the whole network. The percentages shown here are estimates of how much of the network is taken up by each part of the bow tie in the case of the World Wide Web. After Broder *et al.* [84].

10.2 SHORTEST PATHS AND THE SMALL-WORLD EFFECT

One of the most remarkable and widely discussed of network phenomena is the *small-world effect*, the finding that in many networks—perhaps most—the typical distances between pairs of nodes are surprisingly short. In Section 4.6 we discussed Stanley Milgram’s 1967 letter-passing experiment, in which people were asked to get a letter from an initial holder to a distant target person by passing it from acquaintance to acquaintance through the social network. The letters that made it to the target did so in a remarkably small number of steps, around six on average. Milgram’s experiment is a beautiful and convincing demonstration of the small-world effect, though also a rather limited one, given the constraints of the experimental setup. But with the very complete network data we have for many networks these days it is now possible to measure directly the path lengths between nodes and verify the small-world effect explicitly.

Consider an undirected network and, as in Section 7.1.6, let us define d_{ij} to be the length of the shortest path through the network between nodes i and j ,

often simply called the distance between i and j . Then, following Eq. (7.20), the mean distance ℓ_i between node i and every other node is

$$\ell_i = \frac{1}{n} \sum_j d_{ij}. \quad (10.1)$$

Then we define the mean distance ℓ between nodes for the network as a whole as the average of this quantity over all nodes:

$$\ell = \frac{1}{n} \sum_i \ell_i = \frac{1}{n^2} \sum_{ij} d_{ij}. \quad (10.2)$$

In mathematical terms, the small-world effect is the hypothesis that this mean distance is small, in a sense that we will define shortly.

There is a slight catch here, in that Eq. (10.2) does not work for networks with more than one component, because d_{ij} is not well defined for nodes that lie in different components. (Conventionally we sometimes say that such nodes have $d_{ij} = \infty$, but this doesn't help with Eq. (10.2)—it would make ℓ infinite for any network with more than one component.) The most common way around this problem is to change the definition of ℓ to an average only over paths within components. Let \mathcal{C}_m denote the components of a network, with $m = 1, 2, \dots$, and then let us define

$$\ell = \frac{\sum_m \sum_{ij \in \mathcal{C}_m} d_{ij}}{\sum_m n_m^2}, \quad (10.3)$$

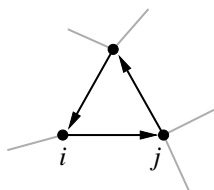
where n_m is the number of nodes in component \mathcal{C}_m . This measure is now finite for all networks, although it is no longer equal to a simple average over the values of ℓ_i for each node.

In Table 10.1 we list the value of ℓ , defined according to Eq. (10.3), for each of the networks in the table, and we see that indeed it takes quite small values, always less than 20 and usually less than 10, even though some of the networks have millions of nodes.

One can well imagine that the small-world effect could have substantial implications for networked systems. Suppose a rumor or a disease is spreading over a social network for instance. Clearly it will reach people much faster if it is only six steps from any person to any other than if it is a hundred, or a million. Similarly, the time it takes to transmit data from one computer to another on the Internet depends on how many steps or “hops” data packets make as they traverse the network. A network in which the typical number of hops is only ten or twenty will perform much better than one in which it is a hundred or more. (While this point was not articulated by the original designers of the

Internet in the 1960s, they must have had some inkling of the idea, to believe that a network like the Internet could be built and made to work.)

Once one looks more deeply into the mathematics of networks, which we will do in later chapters, one discovers that in fact the small-world effect is not so surprising after all. As we will see in Section 11.7, mathematical models of networks suggest that typical path lengths in networks should scale roughly as $\log n$ with the number n of network nodes, and should therefore tend to remain small even for large networks, since the logarithm is a slowly growing function of its argument. Figure 11.7 on page 364, for example, shows measured values of ℓ plotted against $\log n$ for a set of Facebook social networks, and indeed it does appear that ℓ goes roughly as $\log n$ in this case.



The shortest path from i to j in this network has length 1, but the shortest path from j to i has length 2.

One can ask about path lengths in directed networks as well, although the situation is more complicated. Since in general the path from node i to node j is different in a directed network from the path from j to i , the two paths may have different lengths, and both need to be included when calculating the average length ℓ . It's also possible for there to be no path in one or both directions between two nodes, depending on the component structure. As before we can get around the problems of unconnected nodes by defining our averages over only node pairs that are actually connected by a path. Values calculated in this way are given for the directed networks in Table 10.1.

One can also examine the diameter of a network, which, as described in Section 6.11.1, is the length of the longest finite distance between any pair of nodes in the network. The diameter is usually found to be relatively small as well and calculations using network models suggest that it should scale logarithmically with n , as the average distance does. But the diameter is, in general, a less useful measure of real-world network behavior than mean distance because it really only measures the distance between one specific pair of nodes at the extreme end of the distribution of distances. Moreover, the diameter of a network could be affected substantially by a small change to only a single node or a few nodes, which makes it a poor indicator of the typical behavior of the network as a whole.

Nonetheless, there are cases where the diameter is of interest. In Section 10.4 we discuss so-called scale-free networks, i.e., networks with power-law degree distributions. Such networks are believed to have an unusual structure consisting of a central "core" that contains most of the nodes and has a mean distance between node pairs that scales only as $\log \log n$ with network size, and not as $\log n$, making the mean distance for the whole network scale as $\log \log n$ also. Outside of this core there are longer "streamers" or "tendrils" of nodes attached to the core like hair, which have typical length of order $\log n$, making the *diameter* of the network of order $\log n$ [103,114]. This sort of behavior could

be detected by measuring separately the mean path length and diameter of networks of various sizes to confirm that they have different functional forms.²

Another interesting twist on the small-world effect was discussed by Milgram in his original paper on the problem. He noticed, in the course of his letter-passing experiments, that most of the letters destined for a given target person passed through just one or two acquaintances of the target. Thus, it appeared, most people who knew the target person knew him through these one or two people. This idea, that one or two of your acquaintances are especially well connected and responsible for most of the connection between you and the rest of the world, has been dubbed *funneling*, and it too is something we can test against complete networks with the copious data available to us today. If, for instance, we focus on shortest paths between nodes, as we have been doing in this section, then we can measure what fraction of the shortest paths between a node i and every node reachable from it go through each of i 's neighbors in the network. For many networks, this measurement does reveal a funneling effect. For instance, in the coauthorship network of physicists from Table 10.1 we find that, for physicists having five or more collaborators, 48% of shortest paths go through a single neighbor of the average node, the remaining 52% being distributed over the other four or more neighbors. A similar result is seen in the Internet: among nodes having degree five or greater in a May 2005 snapshot of Internet structure at the autonomous system level, an average of 49% of shortest paths go through a single neighbor of the average node. It is tempting to draw conclusions about the routing of Internet packets from this latter result—perhaps that the network will tend to overload a small number of well-connected nodes rather than distributing the load more evenly—but it is worth noting that modern Internet routers incorporate routing algorithms designed specifically to avoid this kind of overload, so simple funneling statistics may not reflect actual traffic patterns very closely.

Milgram referred to these people as “sociometric superstars.” We discussed them previously in Section 4.6.

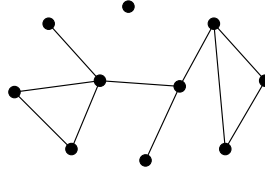
10.3 DEGREE DISTRIBUTIONS

In this section, we look at one of the most fundamental of network properties, the degree distribution. This distribution will come up time and again throughout this book as a defining characteristic of network structure.

Consider an undirected network. As described in Section 6.10, the degree of a node in an undirected network is the number of edges attached to it. The

²It's worth noting, however, that behavior of the form $\log \log n$ is very difficult to confirm in real-world data because $\log \log n$ is a very slowly varying function of n .

degree of an individual in a friendship network, for instance, is just the number of friends they have. Now let us define p_k to be the fraction of nodes that have degree k . For example, consider this network:



It has $n = 10$ nodes, of which one has degree 0, two have degree 1, four have degree 2, two have degree 3, and one has degree 4. Thus the values of p_k are

$$p_0 = \frac{1}{10}, \quad p_1 = \frac{2}{10}, \quad p_2 = \frac{4}{10}, \quad p_3 = \frac{2}{10}, \quad p_4 = \frac{1}{10}, \quad (10.4)$$

and $p_k = 0$ for all $k > 4$. The quantities p_k represent the *degree distribution* of the network. They tell us the frequency with which nodes of different degrees appear in the network.

The value p_k can also be thought of as a probability. It is the probability that a randomly chosen node in the network has degree k . This will be a useful viewpoint when we study theoretical models of networks in Chapters 11 to 13.

Sometimes, rather than the *fraction* of nodes with a given degree, we will want the total number of such nodes. This is easily calculated from the degree distribution simply by multiplying by n . That is, the number of nodes with degree k is np_k , where as usual n is the overall number of nodes in the network.

Another construct containing essentially the same information as the degree distribution is the *degree sequence*, which is the set $\{k_1, k_2, k_3, \dots\}$ of degrees of all the nodes. For instance, the degree sequence of the small network above is $\{0, 1, 1, 2, 2, 2, 2, 3, 3, 4\}$. (The degree sequence need not necessarily be given in ascending order of degrees as here—all orders are equivalent and contain the same information.)

It is probably obvious, but bears saying anyway, that a knowledge of the degree distribution (or degree sequence) does not, in most cases, tell us the complete structure of a network. For most choices of node degrees there is more than one network with those degrees. These two networks, for instance, are different but have the same degrees:



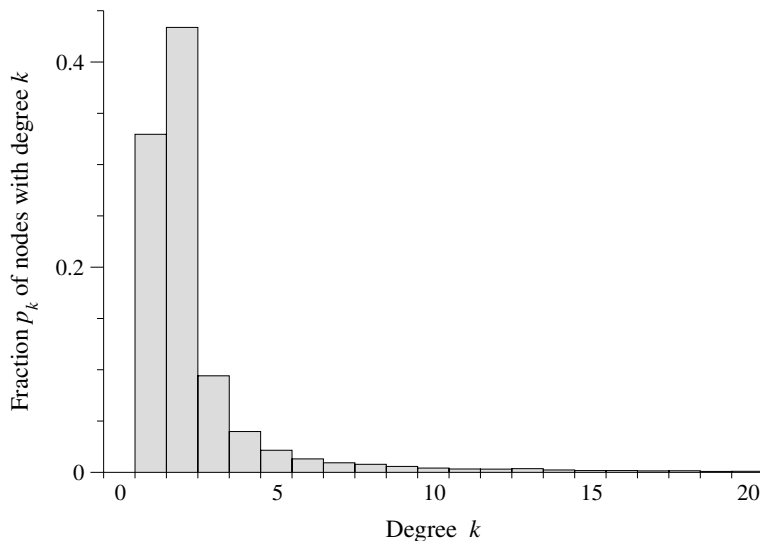


Figure 10.3: The degree distribution of the Internet. A histogram of the degree distribution of the nodes of the Internet at the level of autonomous systems.

Thus we cannot tell the complete structure of a network from its degree distribution alone. The degrees give us important information about a network, but they don't give us complete information.

It is often illuminating to make a plot of the degree distribution of a large network as a function of k . Figure 10.3 shows a plot of p_k for the Internet at the level of autonomous systems. The figure reveals something interesting: most of the nodes in the network have low degree—one or two or three—but there is a significant “tail” to the distribution, corresponding to nodes with substantially higher degree. The plot cuts off at degree 20, but in fact the tail goes much further than this. The highest degree node in the network has degree 2407. Since there are, for this particular data set, a total of 19 956 nodes in the network, that means that the most highly connected node is connected to about 12% of all other nodes in the network. We call such a well-connected node a *hub*.³ Hubs will play an important role in the developments of the

For the Internet there are no nodes of degree zero, since a node is not considered part of the Internet unless it is connected to at least one other.

³We used the word hub in a different and more technical sense in Section 7.1.5 to describe nodes in directed networks that point to many “authorities.” Both senses are common in the networks literature, and in many cases the reader must deduce from the context which is being used. In this book we will mostly use the word in the less technical sense introduced here, of a node with unusually high degree. When we use it in the other sense of Section 7.1.5 we will say so explicitly.

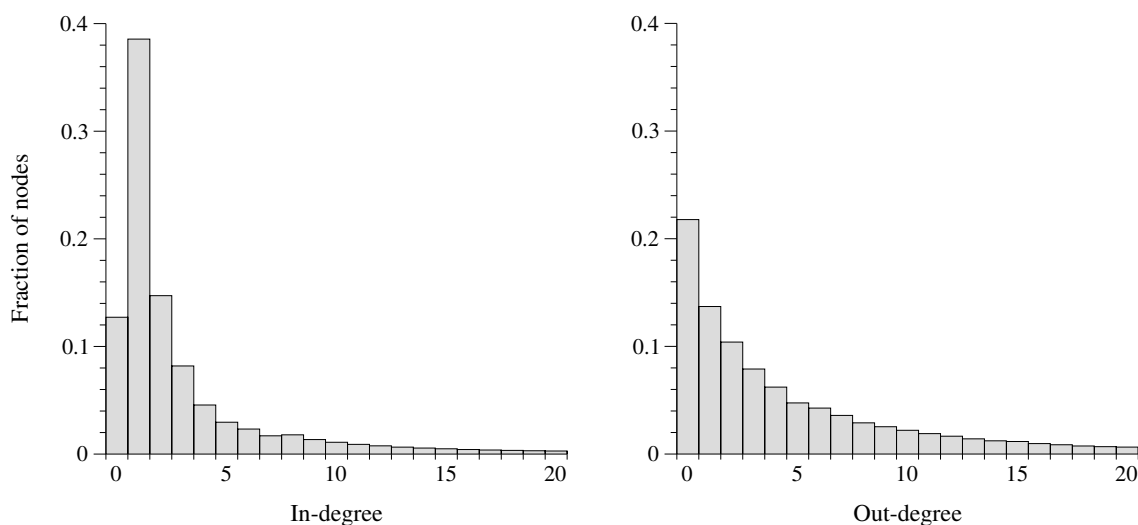


Figure 10.4: The degree distributions of the World Wide Web. Histograms of the distributions of in- and out-degrees of pages on the World Wide Web. Data are from the study by Broder *et al.* [84].

following chapters.

In fact, it turns out that almost all real-world networks have degree distributions with a tail of high-degree hubs like this. In the language of statistics we say that the degree distribution is *right-skewed*. Right-skewed degree distributions are discussed further in Section 10.4, and will reappear repeatedly throughout this book.

One can also calculate degree distributions for directed networks. As discussed in Section 6.10, directed networks have two different degrees for each node, the in-degree and the out-degree, which are, respectively, the number of ingoing and outgoing edges at the node of interest. There are, correspondingly, two different degree distributions in a directed network, the in-degree and out-degree distributions, and one can make a plot of either, or both. Figure 10.4, for example, shows the in- and out-degree distributions for the World Wide Web.

If we wish to be more sophisticated, we might say that the true degree distribution of a directed network is really a joint distribution of in- and out-degrees. We can define p_{jk} to be the fraction of nodes that simultaneously have an in-degree j and an out-degree k . This is a two-dimensional distribution that cannot be plotted as a simple histogram, although it could be represented using a two-dimensional density plot or surface plot. By using a joint distribution we can allow for the possibility that the in- and out-degrees of nodes are correlated.

For instance, if nodes with high in-degree also tend to have high out-degree, then we would see this reflected in large values of p_{jk} when both j and k were large. If we only have the separate distributions of in- and out-degree individually, but not the joint distribution, there is no way of telling whether the network contains such correlations.

In practice, the joint in/out degree distribution of directed networks has rarely been measured or studied, so there is relatively little data on it. This is, in some ways, a pity, since many of our theories of directed networks depend on a knowledge of the joint distribution to give accurate answers (see Section 12.11.1), while others make predictions about the joint distribution that we would like to test against empirical data. For the moment, however, this is an area awaiting more thorough exploration.

10.4 POWER LAWS AND SCALE-FREE NETWORKS

Returning to the Internet, another interesting feature of its degree distribution is shown in Fig. 10.5, where we have replotted the histogram of Fig. 10.3 using logarithmic scales. (That is, both axes are logarithmic. We have also made the bins bigger in the histogram to make the effect clearer—they are of width five in Fig. 10.5 where they were only of width one before.) As the figure shows, when viewed in this way the degree distribution roughly follows a straight line. In mathematical terms that means the logarithm of p_k is a linear function of the logarithm of k thus:

$$\ln p_k = -\alpha \ln k + c, \quad (10.5)$$

where α and c are constants. The minus sign here is optional—we could have omitted it—but it is convenient, since the slope of the line in Fig. 10.5 is clearly negative, making α a positive constant equal to minus the slope in the figure. In this case, the slope gives us a value for α of about 2.1.

Taking the exponential of both sides of Eq. (10.5), we can also write the relation between p_k and k in the form

$$p_k = C k^{-\alpha}, \quad (10.6)$$

where $C = e^c$ is another constant. Distributions of this form, varying as a power of k , are called *power laws*. Based on the evidence of Fig. 10.5 we can say that, roughly speaking, the degree distribution of the Internet follows a power law.

This is, in fact, a common pattern seen in quite a few networks. For instance, as shown in Fig. 10.8 on page 323, both the in- and out-degrees of the World Wide Web roughly follow power-law distributions, as do the in-degrees in many citation networks (but not the out-degrees). As we will see later in this

The power-law distribution is also sometimes called the *Pareto distribution* or *Zipf's law*.

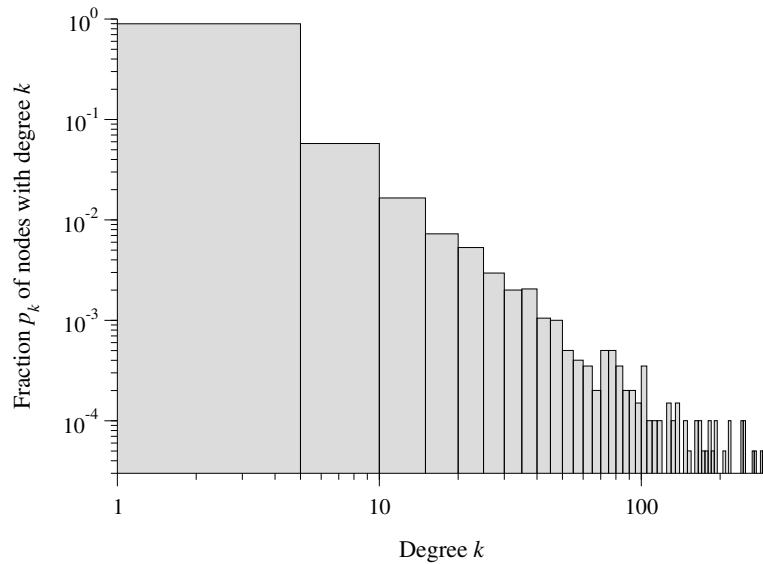


Figure 10.5: The power-law degree distribution of the Internet. Another histogram of the degree distribution of the Internet, plotted this time on logarithmic scales. The approximate straight-line form of the histogram indicates that the degree distribution roughly follows a power law of the form (10.6).

book, networks with power-law degree distributions display some striking and unexpected behaviors.

The constant α is known as the *exponent* of the power law. Values in the range $2 \leq \alpha \leq 3$ are typical, although values a little outside this range are possible and are observed occasionally. Table 10.1 gives the measured values of the exponents for a number of networks that have power-law or approximately power-law degree distributions, and we see that most of them fall in this range. The constant C in Eq. (10.6) is mostly uninteresting, being fixed by the requirement of normalization, as described in Section 10.4.2.

Degree distributions do not usually follow Eq. (10.6) over their entire range. Looking at Fig. 10.3, for example, we can see that the degree distribution falls off at small k . A true power-law distribution is monotonically decreasing over its entire range and hence the degree distribution must in this case deviate from the true power law in the small- k regime. This is typical. A common situation is that the power law is obeyed in the tail of the distribution, for large values of k , but not in the small- k regime. When one says that a particular network has a power-law degree distribution one usually means only that the tail of the

distribution has this form. In some cases, the distribution may also deviate from the power-law form for high k as well. For instance, there may be a cut-off of some type that limits the maximum degree of nodes in the tail.

Networks with power-law degree distributions are sometimes called *scale-free networks*, and we will use this terminology occasionally. Of course, there are many networks that do not have power-law degree distributions, that are not scale-free, but the scale-free ones will be of particular interest to us because of their intriguing properties. Telling the scale-free ones from the non-scale-free is not always easy however. The simplest strategy is to look at a histogram of the degree distribution on a log-log plot, as we did in Fig. 10.5, to see if it follows a straight line. There are, however, a number of problems with this approach and where possible we recommend you use other methods, as we now explain.

10.4.1 DETECTING AND VISUALIZING POWER LAWS

As a tool for visualizing or detecting power-law behavior, a simple histogram like Fig. 10.5 presents some problems. One problem obvious from the figure is that the statistics of the histogram are poor in the tail of the distribution, the large- k region, which is precisely the region in which the power law is normally followed most closely. Each bin of the histogram in this region contains only a few samples, which means that statistical fluctuations in the number of samples from bin to bin are large. This is visible as a “noisy signal” at the right-hand end of Fig. 10.5 that makes it difficult to determine whether the histogram really follows a straight line or not.

There are a number of solutions to this problem. The simplest is to use a histogram with larger bins, so that more samples fall in each bin. In fact, we already did this in going from Fig. 10.3 to Fig. 10.5—we increased the bin width from one to five between the two figures. Larger bins contain more samples and hence give less noise in the tail of the histogram, but at the expense of less detail overall, since the number of bins is correspondingly reduced. Choosing the best bin width can be difficult, in part because different widths may be preferable in different regions of the plot: wide bins might be better for the tail of the distribution where noise is a problem, but we may prefer narrower ones at the left-hand end of the histogram where there are many samples.

We could try to get the best of both worlds by using bins of different sizes in different parts of the histogram. For example, we could use bins of width one for low degrees and switch to width five for higher degrees. In doing this we must be careful to normalize the bins correctly. A bin of width five will on average accrue five times as many samples as a similarly placed bin of width

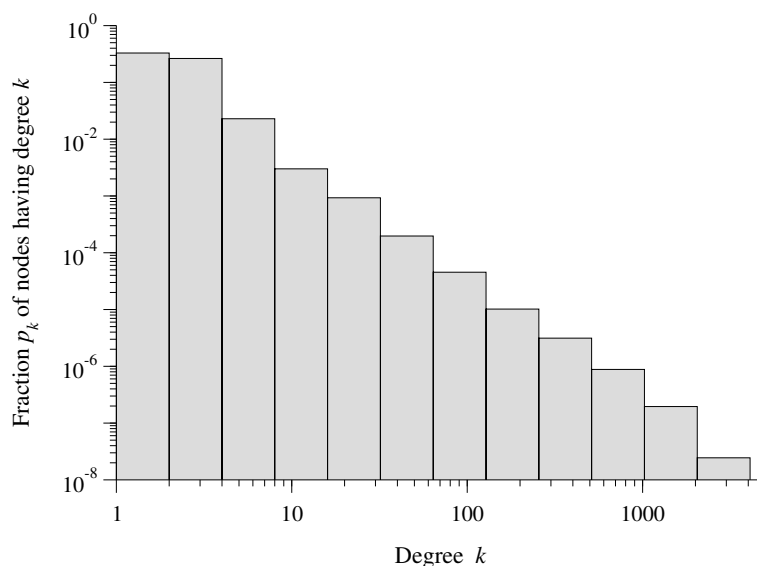


Figure 10.6: Histogram of the degree distribution of the Internet, created using logarithmic binning. In this histogram the widths of the bins are constant on a logarithmic scale, meaning that on a linear scale each bin is wider by a constant factor than the one before it. The counts in the bins are normalized by dividing by bin width to make counts in different bins comparable.

one, so if we wish to plot the counts on the same axes, or otherwise compare them directly, we should divide the number of samples in the larger bin by five. More generally, we should divide sample counts by the width of their bins to make counts in bins of different widths comparable.

We need not restrict ourselves to only two different sizes of bins. We could use larger and larger bins as we go further out in the tail. We can even make every bin a different size, each one a little larger than the one before it. One commonly used version of this idea is called *logarithmic binning*. In this scheme, each bin is made wider than its predecessor by a constant factor a . For instance, if the first bin in a histogram covers the interval $1 \leq k < 2$ (meaning that all nodes of degree 1 fall in this bin) and $a = 2$, then the second would cover the interval $2 \leq k < 4$ (nodes of degrees 2 and 3), the third the interval $4 \leq k < 8$, and so forth. In general the n th bin would cover the interval $a^{n-1} \leq k < a^n$ and have width $a^n - a^{n-1}$. The most common choice for a is $a = 2$, since larger values tend to give bins that are too coarse while smaller ones give bins with non-integer limits.

Figure 10.6 shows the degree distribution of the Internet binned logarithmically in this way. We have been careful to normalize each bin by dividing by its width, as described above. As we can see, the histogram is now much less noisy in the tail and it is considerably easier to see the straight-line behavior of the degree distribution. The figure also reveals a nice property of logarithmically binned histograms, namely that when plotted on logarithmic scales, as here, the bins appear to have equal width. This is, in fact, the principal reason for this particular choice of bins and also the origin of the name “logarithmic binning.”

Note that in a logarithmically binned histogram there is never any bin that contains nodes of degree zero (and even if there were it would not appear in a plot like Fig. 10.6 since there is no zero on a logarithmic scale). If we want to know how many nodes there are of degree zero we will have to measure this number separately.

A different solution to the problem of visualizing a power-law distribution is to construct the *cumulative distribution function*, which is defined by

$$P_k = \sum_{k'=k}^{\infty} p_{k'}. \quad (10.7)$$

In other words, P_k is the fraction of nodes that have degree k or greater.⁴

Suppose the degree distribution p_k follows a power law in its tail. To be precise, let us say that $p_k = C k^{-\alpha}$ for $k \geq k_{\min}$ for some k_{\min} . Then for $k \geq k_{\min}$ we have

$$P_k = C \sum_{k'=k}^{\infty} k'^{-\alpha} \simeq C \int_k^{\infty} k'^{-\alpha} dk' = \frac{C}{\alpha - 1} k^{-(\alpha-1)}, \quad (10.8)$$

where we have approximated the sum by an integral, which is reasonable since the power law is a slowly varying function for large k . (We are also assuming that $\alpha > 1$ so that the integral converges.) Thus, we see that if the distribution p_k follows a power law, then so does the cumulative distribution function P_k , but with an exponent $\alpha - 1$ that is 1 less than the original exponent.

This gives us another way of visualizing a power-law distribution: we plot the cumulative distribution function on logarithmic scales, as we did for the original histogram, and again look for straight-line behavior. We have done

⁴Sometimes one sees the cumulative distribution function defined in the opposite direction as the fraction of nodes with degree k or less, meaning that the sum in Eq. (10.7) would run from zero to k , rather than k to infinity. For our purposes, however, it is important that the function be defined as here, with the sum from k to infinity. Sometimes you may see this function referred to as the *complementary* cumulative distribution function, to distinguish it from the other alternative.

Alternatively, we can think of P_k as the probability that a randomly chosen node has degree k or greater.

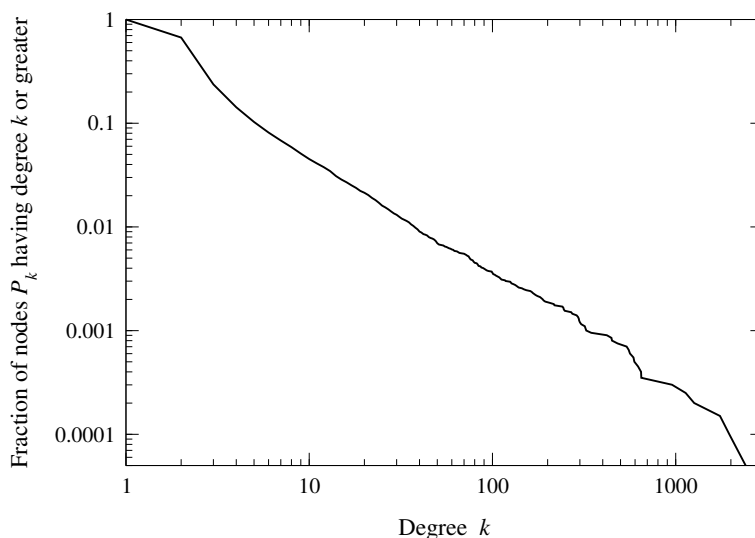


Figure 10.7: Cumulative distribution function for the degrees of nodes on the Internet. For a distribution with a power-law tail, as is approximately the case for the degree distribution of the Internet, the cumulative distribution function, Eq. (10.7), also follows a power law, but with a slope 1 less than that of the original distribution.

this in Fig. 10.7 for the case of the Internet, and the (approximate) straight-line form is clearly visible. Three more examples are shown in Fig. 10.8, for the in- and out-degree distributions of the World Wide Web and for the in-degree distribution of a citation network.

This approach has some advantages. In particular, the calculation of P_k does not require us to bin the values of k as we do with a normal histogram. P_k is perfectly well defined for any value of k and can be plotted just as a normal function. When the bins in a degree histogram contain more than one value of k —i.e., when their width is greater than 1—the binning of the data necessarily throws away some information, eliminating, as it does, the distinction between any two values that fall into the same bin. The cumulative distribution function, on the other hand, preserves all of the information contained in the data, because no bins are involved. The most obvious manifestation of this difference is that the number of points in a plot like Fig. 10.5 or Fig. 10.6 is relatively small, whereas in a cumulative distribution plot like Fig. 10.7 there are as many points along the k (horizontal) axis as there are distinct values of k .

The cumulative distribution function is easy to calculate. We simply sort the

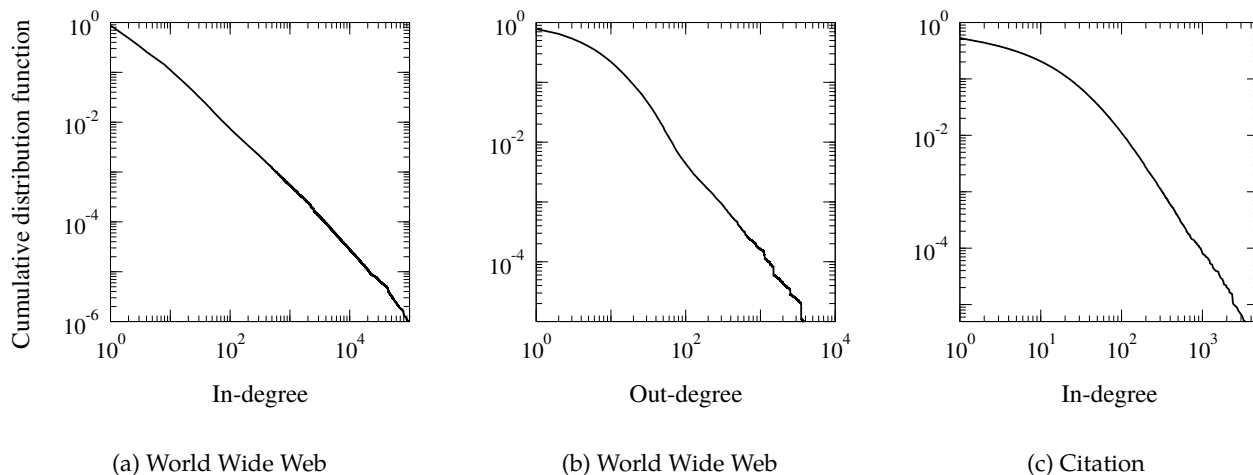


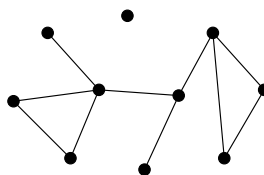
Figure 10.8: Cumulative distribution functions for in- and out-degrees in directed networks. (a) The in-degree distribution of the World Wide Web, from the data of Broder *et al.* [84]. (b) The out-degree distribution for the same web data set. (c) The in-degree distribution of a citation network, from the data of Redner [404]. The distributions follow approximate power-law forms in each case.

degrees of the nodes in descending order and then number them from 1 to n in that order.⁵ These numbers are the so-called *ranks* r_i of the nodes. A plot of r_i/n as a function of degree k_i then gives us our cumulative distribution plot.⁶

For instance, consider again this small example network, which we looked at at the beginning of Section 10.3:

⁵The most time-consuming part of this calculation is the sorting of the degrees. Sorting is a well-studied problem and the fastest general algorithms run in time $O(n \log n)$. Thus, the leading order scaling of this algorithm to calculate the cumulative distribution is $O(n \log n)$. Most computer languages, as well as numerical software such as spreadsheet programs, include built-in functions for sorting numbers, which saves you from having to create your own.

⁶Such plots are also sometimes called *rank-frequency* plots because one of their earliest uses was to detect power-law behavior in the frequency of occurrence of words in natural languages. If the data you are measuring are frequencies, then the cumulative distribution graph is a plot of rank against frequency. More recently such plots have been used to detect power-law behavior in many quantities other than frequencies, but the name “rank-frequency plot” is still often used.



The degrees of the nodes in this case are 0, 1, 1, 2, 2, 2, 2, 3, 3, 4. Sorting these into decreasing order and numbering them, we find values for P_k as follows:

Degree k	Rank r	$P_k = r/n$
4	1	0.1
3	2	0.2
3	3	0.3
2	4	0.4
2	5	0.5
2	6	0.6
2	7	0.7
1	8	0.8
1	9	0.9
0	10	1.0

Then a plot of the last column as a function of the first gives us our cumulative distribution function.

Cumulative distributions do have some disadvantages. One is that they are less easy to interpret than ordinary histograms, since they are only indirectly related to the actual distribution of node degrees. A more serious disadvantage is that successive points on a cumulative distribution plot are correlated—the cumulative distribution function in general only changes a little from one point to the next, so adjacent values are not at all independent. This means that it is not valid, for instance, to extract the exponent of a power-law distribution by fitting the slope of the straight-line portion of a plot like Fig. 10.7 and equating the result with $\alpha - 1$, at least if the fitting is done using standard methods such as least squares that assume independence between the data points.

In fact, it is in general not good practice to evaluate power-law exponents by performing straight-line fits to either cumulative distribution functions or ordinary histograms. Both are known to give biased answers, although for different reasons [111, 209]. Instead, it is usually better to calculate α directly

from the data, using the formula⁷

$$\alpha = 1 + N \left(\sum_i \ln \frac{k_i}{k_{\min} - \frac{1}{2}} \right)^{-1}. \quad (10.9)$$

Here, k_{\min} is the minimum degree for which the power law holds, as before, and N is the number of nodes with degree greater than or equal to k_{\min} . The sum is performed over only those nodes with $k \geq k_{\min}$, and not over all nodes.

We can also calculate the statistical error on α from the formula:

$$\sigma = \sqrt{N} \left(\sum_i \ln \frac{k_i}{k_{\min} - \frac{1}{2}} \right)^{-1} = \frac{\alpha - 1}{\sqrt{N}}. \quad (10.10)$$

For example, applying Eqs. (10.9) and (10.10) to the degree sequence of the Internet from Fig. 10.3 gives an exponent value of $\alpha = 2.11 \pm 0.01$.

The derivation of these formulas, which makes use of maximum likelihood techniques, would take us some way from our primary topic of networks, so we will not go into it here. The interested reader can find a discussion in Ref. [111], along with many other details such as methods for determining the value of k_{\min} and methods for telling whether a particular distribution follows a power law at all.

10.4.2 PROPERTIES OF POWER-LAW DISTRIBUTIONS

Quantities with power-law distributions behave in some surprising ways. We take a few pages here to look at some properties of power-law distributions that will be of use to us later on.

Power laws turn up in a wide variety of places, not just in networks. They are found in the sizes of city populations [32, 486], earthquakes [228], moon craters [344], solar flares [312], computer files [124], and wars [410]; in the frequency of use of words in human languages [164, 486], the frequency of occurrence of personal names in most cultures [480], the numbers of papers scientists write [309], and the number of hits on web pages [5]; in the sales of books, music recordings, and almost every other branded commodity [123, 273]; and in the numbers of species in biological taxa [87, 472]. A review of the data and some mathematical properties of power laws can be found in Ref. [357].

⁷This formula is only an approximation to the full formula for the exponent. The full formula, unfortunately, does not give a closed-form expression for α , making it hard to use. Equation (10.9) works well provided k_{\min} is greater than about 6, which is true for many networks. In cases where it is not, however, the full formula must be used—see Ref. [111].

Here we highlight just a few points that will be relevant for our study of networks.

Normalization: The constant C appearing in the power-law form $p_k = C k^{-\alpha}$, Eq. (10.6), is fixed by the requirement that the degree distribution be normalized. That is, when we add up the total fraction of nodes having all possible degrees $k = 0 \dots \infty$, we must get 1:

$$\sum_{k=0}^{\infty} p_k = 1. \tag{10.11}$$

Even in the best of circumstances our degree distribution cannot obey $p_k = C k^{-\alpha}$ for all k down to zero, because then p_0 would be infinite, which is impossible since it is a probability and must lie between 0 and 1. Let us suppose, however, that the distribution follows the power law for all $k \geq 1$, and that there are no nodes of degree zero so that $p_0 = 0$. Substituting these assumptions in Eq. (10.11) we get $C \sum_{k=1}^{\infty} k^{-\alpha} = 1$, or

$$C = \frac{1}{\sum_{k=1}^{\infty} k^{-\alpha}} = \frac{1}{\zeta(\alpha)}, \tag{10.12}$$

where $\zeta(\alpha) = \sum_{k=1}^{\infty} k^{-\alpha}$ is the Riemann zeta function. Thus the correctly normalized power-law distribution is

$$p_k = \frac{k^{-\alpha}}{\zeta(\alpha)}, \tag{10.13}$$

for $k > 0$, and $p_0 = 0$.

This is a reasonable starting point for mathematical models of scale-free networks—we will use it in Chapter 12, for example—but it’s not a very good representation of most real-world networks, which deviate from pure power-law behavior for small k as described in Section 10.4 and seen in Fig. 10.3. In that case, the normalization constant will take some other value dependent on the particular shape of the distribution, but it is still fixed nonetheless by the normalization condition, Eq. (10.11). For some calculations we will be interested only in the tail of the distribution where the power-law behavior holds and can discard the rest of the data. In such cases, we can normalize over only the tail, starting from the minimum value k_{\min} for which the power law holds. This gives

$$p_k = \frac{k^{-\alpha}}{\sum_{k=k_{\min}}^{\infty} k^{-\alpha}} = \frac{k^{-\alpha}}{\zeta(\alpha, k_{\min})}, \tag{10.14}$$

where $\zeta(\alpha, k_{\min}) = \sum_{k=k_{\min}}^{\infty} k^{-\alpha}$ is the so-called generalized or incomplete zeta function.

Alternatively, we could observe, as we did for Eq. (10.8), that in the tail of the distribution the sum over k is well approximated by an integral, so the normalization constant can be written

$$C \simeq \frac{1}{\int_{k_{\min}}^{\infty} k^{-\alpha} dk} = (\alpha - 1)k_{\min}^{\alpha-1}, \quad (10.15)$$

or

$$p_k \simeq \frac{\alpha - 1}{k_{\min}} \left(\frac{k}{k_{\min}} \right)^{-\alpha}. \quad (10.16)$$

In the same approximation the cumulative distribution function, Eq. (10.8), is given by the simple expression

$$P_k = \left(\frac{k}{k_{\min}} \right)^{-(\alpha-1)}. \quad (10.17)$$

Moments: Of particular interest to us will be the moments of the degree distribution. The first moment of a distribution is its mean:

$$\langle k \rangle = \sum_{k=0}^{\infty} k p_k. \quad (10.18)$$

The second moment is the mean square:

$$\langle k^2 \rangle = \sum_{k=0}^{\infty} k^2 p_k. \quad (10.19)$$

And the m th moment is

$$\langle k^m \rangle = \sum_{k=0}^{\infty} k^m p_k. \quad (10.20)$$

Suppose we have a degree distribution p_k that follows a power law $p_k = Ck^{-\alpha}$ for $k \geq k_{\min}$, in the manner of the Internet or the World Wide Web. Then

$$\langle k^m \rangle = \sum_{k=0}^{k_{\min}-1} k^m p_k + C \sum_{k=k_{\min}}^{\infty} k^{m-\alpha}. \quad (10.21)$$

Since the power law is a slowly varying function of k for large k , we can again approximate the second sum by an integral:

$$\begin{aligned} \langle k^m \rangle &\simeq \sum_{k=0}^{k_{\min}-1} k^m p_k + C \int_{k_{\min}}^{\infty} k^{m-\alpha} dk \\ &= \sum_{k=0}^{k_{\min}-1} k^m p_k + \frac{C}{m - \alpha + 1} \left[k^{m-\alpha+1} \right]_{k_{\min}}^{\infty}. \end{aligned} \quad (10.22)$$

The first term here is some finite number whose value depends on the particular (non-power-law) form of the degree distribution for small k . The second term depends on the values of m and α . If $m - \alpha + 1 < 0$, then the bracket has a finite value, and $\langle k^m \rangle$ is well-defined. But if $m - \alpha + 1 \geq 0$, then the bracket diverges and with it the value of $\langle k^m \rangle$. Thus, the m th moment of the degree distribution is finite if and only if $\alpha > m + 1$. Put another way, for a given value of α all moments will diverge for which $m \geq \alpha - 1$.

Of special concern to us will be the second moment $\langle k^2 \rangle$, which arises in many calculations to do with networks. The results above tell us that the second moment is finite if and only if $\alpha > 3$. As discussed in Section 10.4, however, most real-world networks with power-law degree distributions have values of α in the range $2 \leq \alpha \leq 3$, which means that the second moment should diverge, an observation that has a number of remarkable implications for the properties of scale-free networks, some of which we will explore in coming chapters. Note that this applies even for networks where the power law only holds in the tail of the distribution—the distribution does not have to follow a power law everywhere for the second moment to diverge.

This description, however, is slightly misleading. In any real network all the moments of the degree distribution will actually be finite. We can always calculate the m th moment directly from the degree sequence thus:

$$\langle k^m \rangle = \frac{1}{n} \sum_{i=1}^n k_i^m, \quad (10.23)$$

where k_i is the degree of node i as previously. And since all the k_i are finite, so must this sum be. When we say that the m th moment is infinite, what we really mean is that it would diverge in the limit $n \rightarrow \infty$ of an arbitrarily large network. Even for finite-sized networks, however, the values of the moments, while not infinite, can still become very large, and this alone is enough to produce interesting behaviors. For the Internet data we used in Figs. 10.3 and 10.5, for example, the second moment has the value $\langle k^2 \rangle = 1159$, which can in practice be treated as infinite for many purposes. We will see a number of consequences of such large moment values in later chapters.

Top-heavy distributions: An interesting quantity is the fraction of edges in a network that are connected to the nodes with the highest degrees. For a pure power-law degree distribution, it can be shown [357] that the fraction W of ends of edges attached to the fraction P of highest-degree nodes is

$$W = P^{(\alpha-2)/(\alpha-1)}. \quad (10.24)$$

A set of curves of W against P is shown in Fig. 10.9 for various values of α . Curves of this kind are called *Lorenz curves*, after Max Lorenz, who first studied

The second moment will crop up, for instance, in our calculations of the mean degree of neighbors in Section 12.2, network robustness in Section 15.2.1, and epidemiological processes in Section 16.3.2, as well as a number of other places.

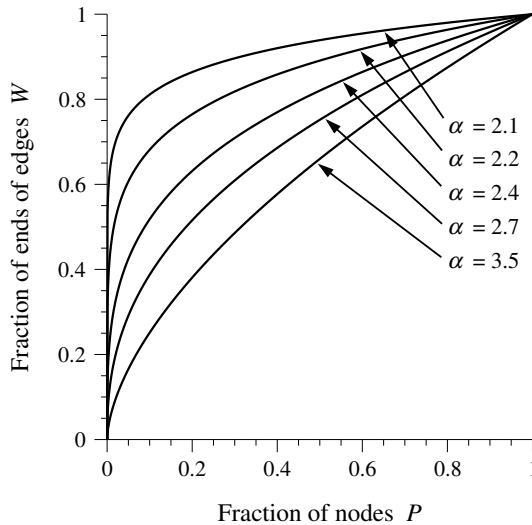


Figure 10.9: Lorenz curves for scale-free networks. The curves show the fraction W of the total number of ends of edges in a scale-free network that are attached to the fraction P of nodes with the highest degrees, for various values of the power-law exponent α .

them around the beginning of the twentieth century [308]. As the figure shows, the curves are concave downward for all values of α , and for values only a little above 2 they have a very steep initial increase, meaning that a large fraction of the edges are connected to a small fraction of the highest degree nodes.

Thus, for example, the in-degree distribution of the World Wide Web follows a power law with exponent around $\alpha = 2.2$. Equation (10.24) with $P = \frac{1}{2}$ then tells us that (assuming a perfect power-law form) we would expect about $W = 0.89$ or 89% of all hyperlinks to go to pages in the top half of the degree distribution, while the bottom half gets a mere 11%. Conversely, if we set $W = \frac{1}{2}$ in Eq. (10.24) we get $P = 0.015$, implying that 50% of all the links go to less than 2% of the “richest” nodes. Thus the degree distribution is “top-heavy,” a large fraction of the “wealth”—meaning incoming hyperlinks in this case—falling to a small fraction of the nodes.⁸

This calculation assumes a degree distribution that follows a perfect power law, whereas in reality, as we have seen, degree distributions usually only follow

⁸Similar results apply to other quantities with power-law distributions, including actual monetary wealth—you may have seen statistics about the large fraction of wealth falling in the hands of the richest people in the population [357].

a power law in their high-degree tail. The basic principle still holds, however, and even if we cannot write an exact formula like Eq. (10.24) for a particular network we can easily evaluate W as a function of P directly from degree data. For the real degree distribution of the Web⁹ we find that 50% of the incoming hyperlinks point to just 1.1% of the richest nodes (so Eq. (10.24) was not far off in this case). Similarly, in scientific citation networks 8.3% of the highest cited papers get 50% of all the citations,¹⁰ and on the Internet just 3.3% of the most highly connected nodes have 50% of the connections.¹¹

In the remaining chapters of this book we will see many examples of networks with power-law degree distributions, and we will make use of the results given here to develop an understanding of their behavior.

10.5 DISTRIBUTIONS OF OTHER CENTRALITY MEASURES

As discussed in Chapter 7, node degree is just one of a variety of centrality measures for nodes in networks. Other measures include eigenvector centrality and its variants (Sections 7.1.2 to 7.1.5), closeness centrality (Section 7.1.6), and betweenness centrality (Section 7.1.7). The distributions of these other measures, while of lesser importance in the study of networks than the degree distribution, are nonetheless of interest.

Eigenvector centrality can be thought of as an extended form of degree centrality, which takes into account not only how many neighbors a node has but also how central those neighbors themselves are (Section 7.1.2). Given its similarity to degree centrality, it is perhaps not surprising to learn that eigenvector centrality, like degree, often has a right-skewed distribution. The left panel of Fig. 10.10 shows the cumulative distribution of eigenvector centralities for the nodes of the Internet, using again the autonomous system level data that we used in Section 10.3. As the figure shows, the tail of the distribution approximately follows a power law, i.e., a straight line on the logarithmic scales used in the plot. Similar power-law behavior is also seen in eigenvector centralities for networks like the World Wide Web and citation networks, while some other networks show right-skewed but non-power-law distributions.

Betweenness centrality (Section 7.1.7) also tends to have right-skewed distributions in most networks. The right panel of Fig. 10.10 shows the cumulative distribution of betweenness for the nodes of the Internet and, as we can see,

⁹Using the data of Broder *et al.* [84].

¹⁰Using the data of Redner [404].

¹¹For the AS-level data of Fig. 10.3.

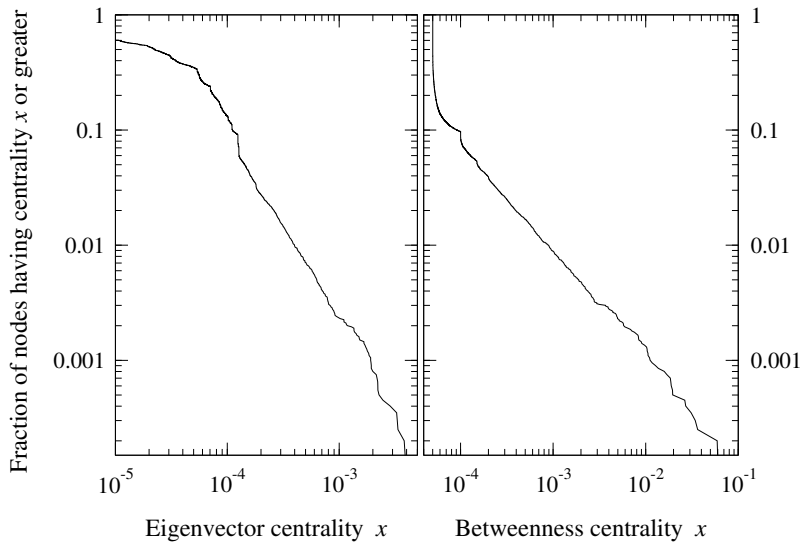


Figure 10.10: Cumulative distribution functions for centralities of nodes on the Internet. Both eigenvector centrality and betweenness appear to roughly follow a power law, at least in the tail of the distribution.

this distribution is again roughly power-law in form. Again there are some other networks that also have power-law betweenness distributions and others still that have skewed but non-power-law distributions.

An exception to this pattern is the closeness centrality (Section 7.1.6), which is the reciprocal of the mean shortest-path distance from a node to all other reachable nodes. The values of the mean distance typically have quite a small range—they are bounded above by the diameter of the network, which, as discussed in Section 10.2, is typically of order $\log n$, and bounded below¹² by 1. This means in practice that the closeness centrality cannot have a broad distribution or a long tail. In Fig. 10.11, for instance, we show a histogram of closeness centrality values for our snapshot of the Internet, and the distribution spans less than an order of magnitude from its minimum value of 0.137 to a maximum of 0.434. There is no long tail to the distribution and it has quite a complicated form with several peaks and dips.

¹²Technically, the lower bound is slightly less than 1. The mean distance from i to all other nodes is $\ell_i = (1/n) \sum_j d_{ij}$ where d_{ij} is the shortest distance between i and j . Noting that $d_{ii} = 0$ and all other $d_{ij} \geq 1$, we then have $\ell_i \geq (n-1)/n$.

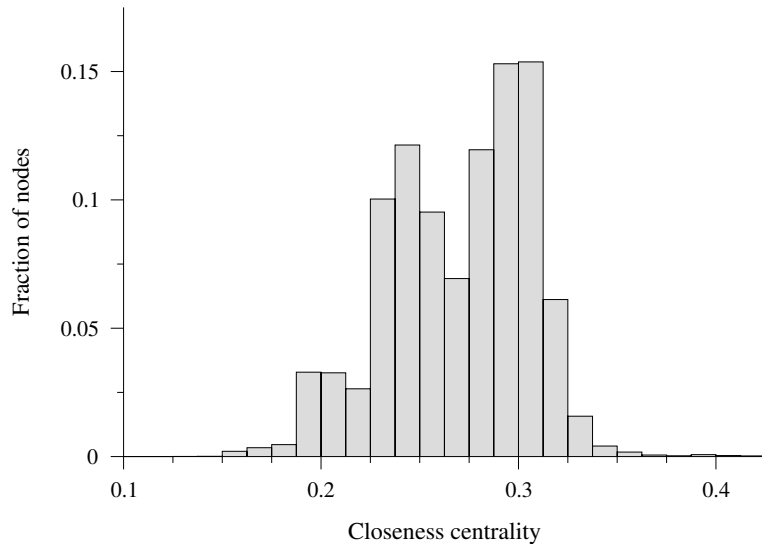


Figure 10.11: Histogram of closeness centralities of nodes on the Internet. Unlike Fig. 10.10 this is a normal non-cumulative histogram showing the actual distribution of closeness centralities. This distribution does not follow a power law.

10.6 CLUSTERING COEFFICIENTS

In Section 7.3 we introduced the clustering coefficient, which is the average probability that two neighbors of the same node are themselves neighbors. The clustering coefficient quantifies the density of triangles in a network and is of interest because in many cases it is found to have values sharply different from what one would expect on the basis of chance. To see what we mean by this, look again at Table 10.1 on page 305, which gives measured values of the clustering coefficient for a variety of networks (in the column denoted C , which gives values for the coefficient defined in Eq. (7.28)). Most of the values are on the order of tens of percent—there is typically a probability between about 10% and 60% that two neighbors of a node will be neighbors themselves.

As we will see in Section 12.3, if we consider a network with a given degree distribution in which connections between nodes are made at random, the clustering coefficient takes the value

$$C = \frac{1}{n} \frac{[\langle k^2 \rangle - \langle k \rangle]^2}{\langle k \rangle^3}, \quad (10.25)$$

where $\langle k \rangle$ and $\langle k^2 \rangle$ are the mean degree and mean-square degree in the net-

work, respectively. Assuming that $\langle k \rangle$ and $\langle k^2 \rangle$ have fixed, well-defined values, this quantity becomes small as $n \rightarrow \infty$ and hence we expect the clustering coefficient to be very small in large networks. This makes the values in Table 10.1, which are relatively large, quite surprising, and indeed many of them turn out to be much larger than the estimate given by Eq. (10.25). For instance, the collaboration network of physicists is measured to have a clustering coefficient of 0.45, but plugging the appropriate values for n , $\langle k \rangle$, and $\langle k^2 \rangle$ into Eq. (10.25) gives $C = 0.0023$. Thus, the measured value is more than a hundred times greater than the value we would expect if physicists chose their collaborators at random.

This large difference is likely indicative of real social effects at work—physicists apparently do not choose their collaborators at random, and moreover choose them in a way that gives rise to a high density of triangles and hence a high value of C . There are a number of reasons why a real collaboration network might contain more triangles than one would expect by chance. One possibility is that people might introduce pairs of their collaborators to each other and those pairs might then go on to collaborate themselves. This is an example of the process sociologists call *triadic closure*: an “open” triad of nodes (i.e., a triad in which one node is linked to the other two, but the third possible edge is missing) is “closed” by the addition of the last edge, forming a triangle.

One can study triadic closure processes directly if one has data on the evolution of a network over time. The network of physics collaborators discussed here was analyzed in this way in Ref. [346], where it was shown that pairs of individuals who have not previously collaborated, but who have another mutual collaborator, are enormously more likely to collaborate in future than pairs who do not—a factor of 45 times as likely in that particular study. Furthermore, this factor increases sharply as the number of mutual collaborators increases, to more than 100 for pairs with two mutual collaborators and almost 150 for pairs with three.

However, it is not always the case that the measured clustering coefficient greatly exceeds the expected value given by Eq. (10.25). Take the example of the Internet. For the Internet at the level of autonomous systems—the same data set we examined in previous sections—the measured clustering coefficient is just 0.012. The expected value if connections were made at random, evaluated from Eq. (10.25), is 0.84. (The large value arises because, as discussed in Section 10.4, the Internet has a highly right-skewed degree distribution, which makes $\langle k^2 \rangle$ large.) Clearly, in this case the clustering is far *less* than one would expect on the basis of chance, suggesting that in the Internet there are forces at work that

The temporal evolution of networks was discussed in Section 6.7.

shy away from the creation of triangles.¹³

In some other networks, such as food webs or the World Wide Web, clustering is neither higher nor lower than expected, taking values roughly comparable with those given by Eq. (10.25). It is not well understood why clustering coefficients take such different values in different types of network, although one theory is that it may be connected with the formation of groups or communities in networks [367].

The clustering coefficient measures the density of triangles in a network. There is no reason, however, for us to limit ourselves to studying only triangles. We can also look at the densities of other small groups of nodes, or *motifs*, as they are often called. One can define coefficients similar to the clustering coefficient to measure the densities of different motifs, although more often one simply counts the numbers of the motifs of interest in a network. And, as with triangles, one can compare the results with the values one would expect to find if connections in the network are made at random. In general, one can find counts that are higher, lower, or about the same as the expected values, all of which can have implications for the understanding of the networks in question. For example, Milo *et al.* [334] looked at motif counts in genetic regulatory networks and neural networks and found certain motifs that occurred far more often than was expected on the basis of chance. They conjectured that these motifs were playing important functional roles in the networks, the equivalent of circuit elements like filters or pulse generators in electronic circuits, and that their frequent occurrence might be an evolutionary result of their usefulness to the organisms involved.

10.6.1 LOCAL CLUSTERING COEFFICIENT

In Section 7.3.1 we introduced the local clustering coefficient C_i for a node i :

$$C_i = \frac{\text{(number of pairs of neighbors of } i \text{ that are connected)}}{\text{(number of pairs of neighbors of } i\text{)}}, \quad (10.26)$$

¹³It is sometimes claimed that essentially all networks show clustering higher than expected [15, 466], which is at odds with the results given here. There seem to be two reasons for the disagreement. First, these claims are based primarily on comparisons of measured clustering coefficients against values calculated on the Poisson random graph, a simple model network with a Poisson degree distribution, which we study in Chapter 11. Many networks, however, have right-skewed degree distributions which are very far from Poisson, and hence the random graph is a poor model against which to compare and probably gives misleading results. Second, the clustering coefficients in these comparisons are mostly calculated as an average of the local clustering, following Eq. (7.31). On networks with highly skewed degree distributions this definition can give very different results from the definition, Eq. (7.28), used in our calculations. Usually Eq. (7.31) gives larger numbers than Eq. (7.28), which could explain the discrepancies in the findings.

which is the fraction of pairs of neighbors of node i that are themselves neighbors. If we calculate the local clustering coefficient for all nodes in a network, in many cases an interesting pattern emerges: we find that on average nodes of higher degree tend to have lower local clustering [402, 458].

Figure 10.12, for example, shows the average value of the local clustering coefficient as a function of degree k on the Internet and the decrease with k in this case is clear. It has been conjectured that plots of this type take either the form $C_i \sim k^{-0.75}$ [458] or the form $C_i \sim k^{-1}$ [402]. In this particular case neither of these conjectures matches the data very well, but for some other networks they appear reasonable.

One possible explanation for the decrease in C_i is that, in some networks at least, nodes tend to clump together into groups or communities, with nodes being connected mostly to others within their own group. (See Chapter 14 for a detailed discussion of the phenomenon of community structure.) In a network showing this kind of behavior, nodes that belong to small groups are constrained to have low degree, because they have relatively few fellow group members to connect to, while those in larger groups can have higher degree. (They don't have to have higher degree, but they can.) At the same time, the local clustering coefficient of nodes in small groups will tend to be larger because each group, being mostly detached from the rest of the network, functions roughly as its own small network and, as discussed earlier, smaller networks are expected to have higher clustering (see Eq. (10.25) and the accompanying discussion). In a network with many groups spanning a range of different sizes, therefore, we would expect nodes of lower degree to have higher clustering on average, as in Fig. 10.12.¹⁴

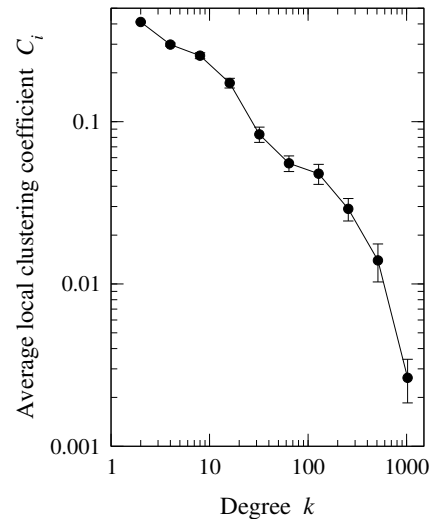


Figure 10.12: Local clustering as a function of degree on the Internet. A plot of the measured mean local clustering coefficient of nodes on the Internet (at the level of autonomous systems) as a function of node degree.

10.7 ASSORTATIVE MIXING

Assortative mixing or homophily is the tendency of nodes to connect to others that are like them in some way. We discussed assortative mixing in Section 7.7 and gave examples such as the high school friendships depicted in Figs. 7.12 and 7.13, where school students tend to associate more with others of the same

¹⁴An alternative proposal is that the behavior of the local clustering coefficient arises through hierarchical structure in a network—that not only are there groups, but that the groups are divided into smaller groups, and those into still smaller ones, and so on. See Refs. [144, 402, 443].

ethnicity or age as themselves.

Of particular interest is assortative mixing by degree, the tendency of nodes to connect others with degrees similar to their own, as discussed in Section 7.7.3. We can also have disassortative mixing by degree, in which nodes connect to others with very different degrees. Both assortative and disassortative mixing can have substantial effects on network structure—see Fig. 7.14 on page 210.

Assortative mixing by degree can be quantified in a number of different ways. One of them is to use the correlation coefficient defined in Eq. (7.64):

$$r = \frac{\sum_{ij}(A_{ij} - k_i k_j / 2m)k_i k_j}{\sum_{ij}(k_i \delta_{ij} - k_i k_j / 2m)k_i k_j}. \quad (10.27)$$

If we are going to calculate the value of this coefficient, however, we should not do it directly from this equation, because the double sum over nodes i and j has a lot of terms (n^2 of them) and is slow to evaluate on a computer. Instead we write

$$r = \frac{S_1 S_e - S_2^2}{S_1 S_3 - S_2^2}, \quad (10.28)$$

with

$$S_e = \sum_{ij} A_{ij} k_i k_j = 2 \sum_{\text{edges } (i,j)} k_i k_j, \quad (10.29)$$

where the second sum is over all distinct (unordered) pairs of nodes (i, j) connected by an edge, and

$$S_1 = \sum_i k_i, \quad S_2 = \sum_i k_i^2, \quad S_3 = \sum_i k_i^3. \quad (10.30)$$

The sum in (10.29) has m terms, where m is the number of edges in the network, and the sums in (10.30) have n terms each. Since $m \ll n^2$ on a typical sparse network, Eq. (10.28) is usually a lot faster to evaluate than Eq. (10.27).

In Table 10.1 we show the values of r for a range of networks and the results reveal an interesting pattern. While none of the values are of very large magnitude—the correlations between degrees are not especially strong—there is a clear tendency for the social networks to have positive r , indicating assortative mixing by degree, while the rest of the networks—technological, information, biological—tend to have negative r , indicating disassortative mixing.

The reasons for this pattern are not known for certain, but it appears that many networks have a tendency to negative values of r because they are simple networks; i.e., they have only single edges between nodes, not multiedges. As shown by Maslov *et al.* [323], networks that have only single edges tend in

the absence of other biases to show disassortative mixing by degree because the number of edges that can fall between high-degree node pairs is limited. Since most networks are represented as simple networks this implies that most should be disassortative, as Table 10.1 indicates.

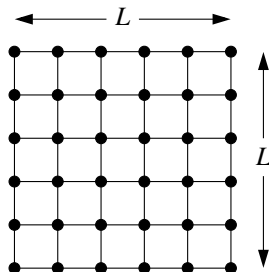
And what about the social networks? One suggestion is that social networks are assortatively mixed because their nodes tend to divide into groups, as discussed in Section 10.6.1. If a network is composed of groups of nodes such that most edges fall within groups, then, as we have said, nodes in small groups tend to have lower degree than nodes in larger groups, simply because nodes in smaller groups have fewer fellow group members to connect to. But since the members of small groups are in groups with other members of the same small groups, it follows that the low-degree nodes will tend to be connected to other low-degree nodes, and similarly for high-degree ones. This simple idea can be turned into a quantitative calculation and indeed it appears that, at least in some circumstances, this mechanism does produce positive values of r [367].

Thus, a possible explanation of the pattern of r -values seen in Table 10.1 is that most networks are naturally disassortative by degree because they are simple networks, while social networks (and perhaps a few others) override this natural bias and become assortative by virtue of their group structure.

EXERCISES

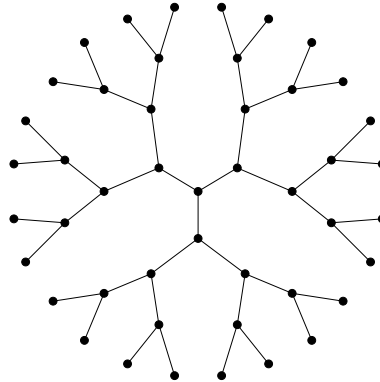
10.1 One can calculate the diameter of certain types of networks exactly.

- a) What is the diameter of a clique?
- b) What is the diameter of a square portion of square lattice, with L edges (or equivalently $L + 1$ nodes) along each side, like this:



What is the diameter of the corresponding hypercubic lattice in d dimensions with L edges along each side? Hence what is the diameter of such a lattice as a function of the number n of nodes?

- c) A Cayley tree is a symmetric regular tree in which each node is connected to the same number k of others, until we get out to the leaves, like this:



(We have $k = 3$ in this picture.) Show that the number of nodes reachable in d steps from the central node is $k(k - 1)^{d-1}$ for $d \geq 1$. Hence find an expression for the diameter of the network in terms of k and the number of nodes n .

- d) Which of the networks in parts (a), (b), and (c) displays the small-world effect, defined as having a diameter that increases as $\log n$ or slower?

10.2 Suppose that a network has a degree distribution that follows the exponential (or geometric) form $p_k = Ca^k$, where C and a are positive constants and $a < 1$.

- Assuming the distribution is properly normalized, find C as a function of a .
- Calculate the fraction P of nodes that have degree k or greater.
- Calculate the fraction W of ends of edges that are attached to nodes of degree k or greater.
- Hence show that the Lorenz curve—the equivalent of Eq. (10.24) for this degree distribution—is given by

$$W = P - \frac{1 - 1/a}{\log a} P \log P.$$

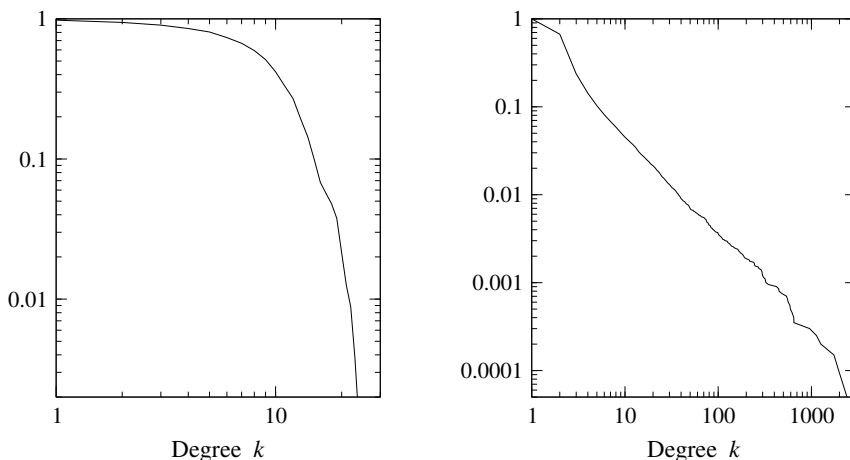
- Show that the value of W is greater than one for some values of P in the range $0 \leq P \leq 1$. What is the meaning of these “unphysical” values?

10.3 A particular network is believed to have a degree distribution that follows a power law for nodes of degree 10 or greater. Among a random sample of nodes in the network, the degrees of the first 20 nodes with degree 10 or greater are:

16	17	10	26	13
14	28	45	10	12
12	10	136	16	25
36	12	14	22	10

Estimate the exponent α of the power law and the error on that estimate using Eqs. (10.9) and (10.10).

10.4 Here are plots of the cumulative distribution function of degrees in two undirected networks:



- One of these networks is approximately scale-free; the other is not. Which is which and how can you tell?
- For the scale-free network give an estimate of the exponent α of the degree distribution.
- If there are m edges in the scale-free network, then there are $2m$ ends of edges. Approximately what fraction of the highest-degree nodes have a half of all the edge-ends?

10.5 Give a derivation of Eq. (10.28) starting from the definition of the correlation coefficient in Eq. (10.27).

10.6 Consider the following simple and rather unrealistic mathematical model of a network. Each of n nodes belongs to one of several groups. The m th group has n_m nodes and each node in that group is connected to others in the group with independent probability $p_m = A(n_m - 1)^{-\beta}$, where A and β are constants, but not to any nodes in other groups. Thus this network takes the form of a set of disjoint clusters or communities.

- Calculate the average degree $\langle k \rangle$ of a node in group m .
- Calculate the average value \bar{C}_m of the local clustering coefficient for nodes in group m .
- Hence show that $\bar{C}_m \propto \langle k \rangle^{-\beta/(1-\beta)}$.
- What value would β have to have for the average value of the local clustering to fall off with increasing degree as $\langle k \rangle^{-3/4}$?

PART III

NETWORK MODELS

CHAPTER 11

RANDOM GRAPHS

An introduction to the most basic of network models, the random graph

IN THE preceding chapters of this book we have looked at how we measure the structure of networks and at mathematical, statistical, and computational methods for making sense of the data we get from our measurements. We have seen, for instance, how to measure the structure of the Internet, and once we have measured it how to determine its degree distribution, its diameter, or the centrality of its nodes. An obvious next question to ask is, “If I know a network has some particular property, such as a particular degree distribution, what effect will that have on the broader behavior of the system?” It turns out that properties like degree distributions can in fact have huge effects on networked systems, which is one of the main reasons we are interested in them. And one of the best ways to understand and get a feel for these effects is to build mathematical models.

The next few chapters of this book are devoted to an examination of some of the most widely used models of network structure, models that mimic the patterns of connections in networks. These models allow us to create artificial networks with known parameters, which are useful for two reasons: first, they give us insight into fundamental questions of structure—why networks look the way they do and how they change with changing parameters—and, second, they provide a foundation on which to build an understanding of processes taking place on networks, such as diseases spreading on social networks or search engines searching the Web.

Networks, 2nd edition. Mark Newman, Oxford University Press (2018). © Mark Newman.
DOI: 10.1093/oso/9780198805090.001.0001

In Section 10.4, for instance, we noted that many networks have degree distributions that roughly follow a power law—these are the so-called scale-free networks. A reasonable question to ask would be how the structure and behavior of such scale-free networks differs from that of their non-scale-free counterparts. One way to address this question would be to create, on a computer for example, two artificial networks, one with a power-law degree distribution and one without, and explore their differences empirically. Better still, one could create a large number of networks in each of the two classes, to see what statistically significant features appear in one class and not in the other. This is precisely the rationale behind random graph models, which are the topic of this chapter and the following one. With random graph models, one creates networks that possess particular properties of interest, such as specified degree distributions, but which are otherwise random.

In Chapter 13 we look at a different class of models, models in which a network grows or evolves according to a specified set of rules. Such models are particularly useful for understanding how network structure arises in the first place. By growing networks according to a variety of different rules, for example, then comparing the results with observed networks, we can get a feel for which growth processes might be at work in real networks. In later chapters of the book, we will also introduce a number of further models as the need arises in our study of network processes.

11.1 RANDOM GRAPHS

A *random graph* is a model network in which the values of certain properties are fixed, but the network is in other respects random. One of the simplest examples of a random graph is the one where we fix only the number of nodes n and the number of edges m . That is, we take n nodes and place m edges among them at random. More precisely, we choose m distinct pairs of nodes uniformly at random from all possible pairs and connect them with an edge. This model is often referred to by its mathematical name $G(n, m)$.

Another entirely equivalent definition of the model is to say that the network is created by choosing uniformly at random among the set of all simple networks with exactly n nodes and m edges. Since there are $\binom{n}{2}$ pairs of nodes between which we could place an edge, there are $\binom{\binom{n}{2}}{m}$ ways of placing the m edges, and we simply choose any one of these with equal probability.

Strictly, in fact, a random graph model is defined not in terms of a single randomly generated network, but as an *ensemble* of networks, i.e., a probability distribution over possible networks. Thus the model $G(n, m)$ is correctly

Recall that a simple network is one with no multi-edges or self-edges—see Section 6.1.

defined as a probability distribution $P(G)$ over all networks G with

$$P(G) = \frac{1}{\binom{\binom{n}{2}}{m}} \quad (11.1)$$

for simple networks with n nodes and m edges and zero otherwise. (We will see more complicated examples of random graph ensembles shortly.)

When one talks about the properties of random graphs one typically means the average properties of the ensemble. For instance, the diameter of a random graph normally means the diameter $\ell(G)$ of an individual network G averaged over the probability distribution of G thus:

$$\langle \ell \rangle = \sum_G P(G) \ell(G). \quad (11.2)$$

This is a useful definition for several reasons. First, it turns out to lend itself well to analytic calculations; many such average properties of random graphs can be calculated exactly, at least in the limit of large network size. Second, it often reflects exactly the thing we want to get at in making our model network in the first place. Very often we are interested in the typical properties of networks. We might want to know, for instance, what the typical diameter is of a network with a given number of edges. Certainly there are special cases of such networks that have particularly large or small diameters, but these don't reflect the typical behavior. If it's typical behavior we are after, then the ensemble average of a property is often a good guide. Third, it can be shown that the distribution of values for many network measures is sharply peaked, becoming concentrated more and more narrowly around the ensemble average as the size of the network becomes large, so that in the large- n limit, essentially all values one is likely to encounter are very close to the mean.

Some average properties of the random graph $G(n, m)$ are straightforward to calculate. Obviously the average number of edges is m , for instance, and the average degree is $2m/n$. Unfortunately, other properties are not so easy to calculate, and most mathematical work has actually been conducted on a slightly different model that is easier to handle. This model is called $G(n, p)$. In $G(n, p)$ we fix not the number but the *probability* of edges between nodes. Again we have n nodes, but now we place an edge between each distinct pair with independent probability p . In this model the number of edges is not fixed. Indeed it is possible that the network could have no edges at all, or it could have edges between every distinct pair of nodes. (For most values of p these are not likely outcomes, but they could happen.)

Again, the technical definition of the model is not in terms of a single network, but in terms of the ensemble, the probability distribution over all

possible networks. Specifically, $G(n, p)$ is the ensemble of simple networks with n nodes in which each network G appears with probability

$$P(G) = p^m(1 - p)^{\binom{n}{2} - m}, \quad (11.3)$$

where m is the number of edges in the network.

$G(n, p)$ was first studied, to this author's knowledge, by Solomonoff and Rapoport [433], but it is most closely associated with the names of Paul Erdős and Alfréd Rényi, who published a celebrated series of papers about the model in the late 1950s and early 1960s [158–160]. If you read scientific papers on this subject, you will sometimes find the model referred to as the “Erdős–Rényi model” or the “Erdős–Rényi random graph” in honor of their contribution. It is also sometimes called the “Poisson random graph” or the “Bernoulli random graph,” names that refer to the distributions of degrees and edges in the model respectively. And sometimes the model is referred to simply as “the” random graph—there are many random graph models, but $G(n, p)$ is the most fundamental and widely studied of them, so if someone is talking about a random graph but doesn't bother to mention which one, they are probably thinking of this one.

The random graph $G(n, p)$ is a very simple model of a network, and there are many features of real networks that it fails to capture. We will see some of them shortly. As is the case in many branches of science, however, we can learn a lot by studying simple models. In this chapter we describe the basic mathematics of $G(n, p)$, focusing particularly on its degree distribution and component structure, which are two of the model's most illuminating characteristics. The techniques we develop in this chapter will also prove useful for some of the more complex models examined later in the book.

11.2 MEAN NUMBER OF EDGES AND MEAN DEGREE

As an example of a very simple calculation for the random graph $G(n, p)$ let us compute the expected number of edges in our model network. We have said that the number of edges m in the model is not fixed, but we can calculate its mean $\langle m \rangle$ easily enough: the average number of edges between a single pair of nodes is p by definition, and the average number between all $\binom{n}{2}$ pairs is simply $\binom{n}{2}$ times this, or

$$\langle m \rangle = \binom{n}{2} p. \quad (11.4)$$

We can use this result to also calculate the mean degree of a node. The mean degree in a network with exactly m edges is $2m/n$ (see Eq. (6.15)) and the mean

degree $\langle k \rangle$ in $G(n, p)$ is the average of this quantity:

$$\langle k \rangle = \left\langle \frac{2m}{n} \right\rangle = \frac{2\langle m \rangle}{n} = \frac{2}{n} \binom{n}{2} p = (n-1)p, \quad (11.5)$$

where we have used Eq. (11.4) and the fact that n is constant. Previously in this book we have denoted the mean degree by c , and we will adopt this convention here also, writing

$$c = (n-1)p. \quad (11.6)$$

In other words, the average number of edges connected to a node is equal to the expected number p between the node and any other node, multiplied by the number $n-1$ of other nodes.

11.3 DEGREE DISTRIBUTION

Only slightly more taxing is the calculation of the degree distribution of $G(n, p)$. A given node in the network is connected with independent probability p to each of the $n-1$ other nodes. Thus the probability of being connected to a particular k other nodes and not to any of the remainder is $p^k(1-p)^{n-1-k}$. There are $\binom{n-1}{k}$ ways to choose those k other nodes, and hence the total probability of being connected to exactly k others is

$$p_k = \binom{n-1}{k} p^k (1-p)^{n-1-k}, \quad (11.7)$$

which is a binomial distribution. In other words, $G(n, p)$ has a binomial degree distribution.

In many cases we are interested in the properties of large networks, so that n can be assumed to be large. Furthermore, as discussed in Section 6.10, almost all real-world networks are sparse, meaning that only a tiny fraction of the $\binom{n}{2}$ possible edges are actually present and the average degree c is much less than n . More formally, a sparse network is one in which the average degree increases slower than n as n becomes large, in which case Eq. (11.6) implies that $p = c/(n-1)$ will become vanishingly small, which allows us to write

$$\begin{aligned} \ln[(1-p)^{n-1-k}] &= (n-1-k) \ln\left(1 - \frac{c}{n-1}\right) \\ &\simeq -(n-1-k) \frac{c}{n-1} \simeq -c, \end{aligned} \quad (11.8)$$

where we have expanded the logarithm as a Taylor series, and the equalities become exact as $n \rightarrow \infty$ with k fixed. Taking exponentials of both sides, we

then find that $(1 - p)^{n-1-k} = e^{-c}$ in the large- n limit. Also for large n we have

$$\binom{n-1}{k} = \frac{(n-1)!}{(n-1-k)!k!} \simeq \frac{(n-1)^k}{k!}, \quad (11.9)$$

and thus Eq. (11.7) becomes

$$p_k = \frac{(n-1)^k}{k!} p^k e^{-c} = \frac{(n-1)^k}{k!} \left(\frac{c}{n-1}\right)^k e^{-c} = e^{-c} \frac{c^k}{k!}, \quad (11.10)$$

in the limit of large n .

Equation (11.10) is the Poisson distribution. In the limit of large n , $G(n, p)$ has a Poisson degree distribution. This is the origin of the name *Poisson random graph* mentioned in Section 11.1, which we will use occasionally to distinguish this model from some of the other random graph models introduced in following chapters, which don't in general have Poisson degree distributions.

11.4 CLUSTERING COEFFICIENT

A very simple quantity to calculate for the Poisson random graph is the clustering coefficient. Recall that the clustering coefficient C is a measure of the transitivity in a network (Section 7.3) and is defined as the probability that two network neighbors of a node are also neighbors of each other. In a random graph the probability that *any* two nodes are neighbors is exactly the same—all such probabilities are equal to $p = c/(n-1)$. Hence

$$C = \frac{c}{n-1}. \quad (11.11)$$

This is one of several respects in which the random graph differs sharply from most real-world networks. Real-world networks often have quite high clustering coefficients—see Table 10.1 on page 305—while Eq. (11.11) tends to zero in the limit $n \rightarrow \infty$ if the network is sparse (meaning that c increases more slowly than n with growing network size). This discrepancy is discussed further in Section 11.8.

11.5 GIANT COMPONENT

Consider the Poisson random graph $G(n, p)$ for $p = 0$, as shown in Fig. 11.1a. In this case there are no edges in the network at all and it is completely disconnected. Each node is an island on its own and the network has n separate components of exactly one node each.

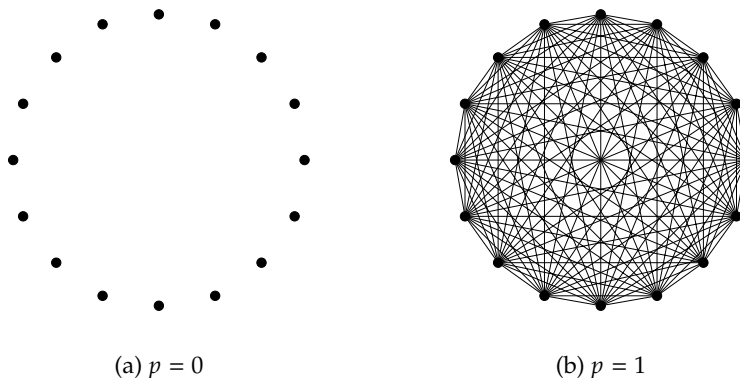


Figure 11.1: Limits of the random graph for $p = 0$ and $p = 1$. These two figures show two random graphs with their nodes arranged in a circle. (a) When $p = 0$ the random graph has no edges, every node is its own component, and the largest component has size 1. (b) When $p = 1$ every possible edge is present, all nodes belong to a single component, and the largest component has size n .

Figure 11.1b shows the opposite limit of $p = 1$, in which every possible edge in the network is present and the network is an n -node clique in the technical sense of the word (Section 7.2.1), meaning that every node is connected directly to every other. In this case, the nodes form a single component that spans the entire network.

The *largest* component in the network has size 1 in the first of these situations ($p = 0$) and size n in the second ($p = 1$). Apart from one of these being much larger than the other, there is an important qualitative difference between these two cases: in the first case the size of the largest component is independent of the number of nodes n in the network, while in the second it is proportional to n , or *extensive* in the scientific jargon. A network component whose size grows in proportion to n we call a *giant component*.

This distinction between the two cases is an important one. In many applications of networks it is crucial that there be a component that fills most of the network. It doesn't necessarily have to fill the entire network, but it should at least fill a large fraction. For instance, in the Internet it is important that there be a path through the network from most computers to most others. If there were not, the network wouldn't be able to perform its intended role of providing computer-to-computer communications for its users. Moreover, as discussed in Section 10.1, most networks do in fact have a large component that fills most

of the network. We can gain some useful insights about what is happening in such networks by considering how the components in our random graph behave.

So let us consider the largest component of our random graph, which has constant size 1 when $p = 0$ and extensive size n when $p = 1$. An interesting question to ask is how the transition between these two extremes occurs if we gradually increase the value of p , starting at 0 and ending up at 1. We might guess, for instance, that the size of the largest component also increases gradually, becoming truly extensive only in the limit where $p = 1$. In reality, however, something much more interesting happens. As we will see, the size of the largest component undergoes a sudden change, or *phase transition*, from constant size to extensive size at one particular value of p . Let us take a look at this transition.

Suppose, for some value of p there is a giant component in the network, meaning that as the size of the network grows, the average size of the largest component grows with it, and hence the largest component occupies a constant fraction of the whole network. We can calculate this fraction exactly in the limit of large network size $n \rightarrow \infty$ as follows.

Let us denote by u the average fraction of nodes in the random graph that do *not* belong to the giant component. Alternatively, we can regard u as the probability that a randomly chosen node in the network does not belong to the giant component. For a node i to not belong to the giant component it must not be connected to the giant component via any other node. If it has even a single edge connecting it to a node in the giant component, then it is itself also in the giant component.

This means that for every other node j in the network either (a) i is not connected to j by an edge, or (b) i is connected to j but j is itself not a member of the giant component. The probability of outcome (a) is simply $1 - p$, the probability of not having an edge between i and j . The probability of outcome (b) is pu , where the factor of p is the probability of having an edge and the factor u is the probability that node j doesn't belong to the giant component.¹ Thus, the complete probability of not being connected to the giant component via node j is $1 - p + pu$.

Then the total probability u of not being connected to the giant component

¹We need to be a little careful here: we really want the probability that j is not connected to the giant component via any node other than node i . However, it turns out that in the limit of large system size this probability is just equal to u . For large n the probability of not being connected to the giant component via any of the $n - 2$ nodes other than i is not significantly smaller than the probability for all $n - 1$ nodes.

via any of the $n - 1$ other nodes in the network is just this quantity to the power of $n - 1$, or

$$u = (1 - p + pu)^{n-1}. \quad (11.12)$$

This gives us a self-consistent equation whose solution tells us the value of u .

We can simplify this equation a bit further as follows. First, let us rearrange it slightly and substitute for p from Eq. (11.6) thus:

$$u = \left[1 - \frac{c}{n-1}(1-u) \right]^{n-1}. \quad (11.13)$$

Now we take logs of both sides:

$$\begin{aligned} \ln u &= (n-1) \ln \left[1 - \frac{c}{n-1}(1-u) \right] \\ &\simeq -(n-1) \frac{c}{n-1}(1-u) = -c(1-u), \end{aligned} \quad (11.14)$$

where the approximate equality becomes exact in the limit of large n . Taking exponentials of both sides again, we then find that

$$u = e^{-c(1-u)}. \quad (11.15)$$

But if u is the fraction of nodes not in the giant component, then the fraction of nodes that are in the giant component is $S = 1 - u$. Eliminating u in favor of S then gives us

$$S = 1 - e^{-cS}. \quad (11.16)$$

This equation, which was first given by Erdős and Rényi in 1959 [158], tells us the size of the giant component as a fraction of the size of the network in the limit of large network size, for any given value of the mean degree c . Unfortunately, though the equation is itself quite simple, it doesn't have a simple solution for S in closed form.² We can, however, get a good feeling for

²One can write a closed-form solution in terms of the *Lambert W-function*, which is defined as the solution $W(z)$ of the equation $W(z)e^{W(z)} = z$. In terms of this function the size of the giant component is

$$S = 1 + \frac{W(-ce^{-c})}{c},$$

where we take the principal branch of the W -function. This expression may have some utility for numerical calculations and series expansions, but it is not widely used. Alternatively, while we cannot write a simple solution for S as a function of c , we can write a solution for c as a function of S . Rearranging Eq. (11.16) for c gives

$$c = -\frac{\ln(1-S)}{S},$$

which can be useful, for instance, for plotting purposes. (We can make a plot of S as a function of c by first making a plot of c as a function of S and then swapping the axes.)

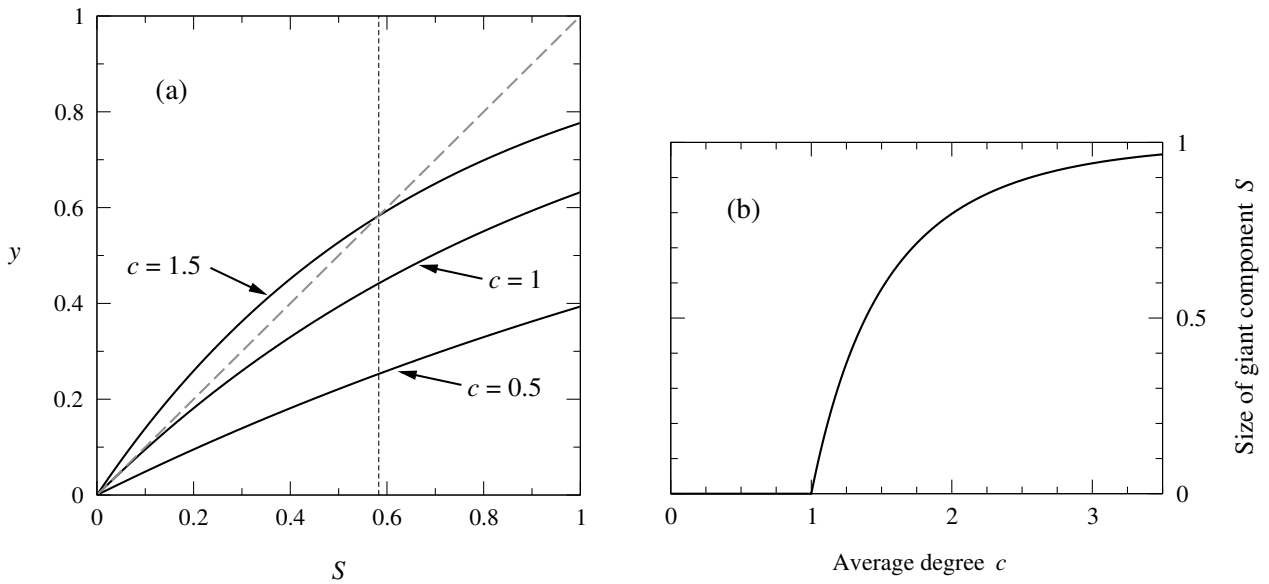


Figure 11.2: Graphical solution for the size of the giant component. (a) The three curves show $y = 1 - e^{-cS}$ for values of c as marked, the diagonal dashed line shows $y = S$, and the intersection gives the solution to Eq. (11.16), $S = 1 - e^{-cS}$. For the bottom curve there is only one intersection, at $S = 0$, so there is no giant component, while for the top curve there is a solution at $S = 0.583\dots$ (vertical dashed line). The middle curve is precisely at the threshold between the regime where a non-trivial solution for S exists and the regime with only the trivial solution $S = 0$. (b) The resulting solution for the size of the giant component as a function of c .

its behavior from a graphical solution. Consider Fig. 11.2a. The three curves show the function $y = 1 - e^{-cS}$ for different values of c . Note that S can take only values from zero to one, so only this part of the curve is shown. The dashed line in the figure is the function $y = S$. Where line and curve cross we have $S = 1 - e^{-cS}$ and the corresponding value of S is a solution to Eq. (11.16).

As the figure shows, depending on the value of c there may be either one solution for S or two. For small c (bottom curve in the figure) there is just one solution at $S = 0$, which implies that there is no giant component in the network. (You can confirm for yourself that $S = 0$ is a solution directly from Eq. (11.16).) On the other hand, if c is large enough (top curve) then there are two solutions, one at $S = 0$ and one at $S > 0$. Only in this regime can there be a giant component.

The transition between the two regimes corresponds to the middle curve in the figure and falls at the point where the gradient of the curve and the gradient

of the dashed line match at $S = 0$. That is, the transition takes place when

$$\frac{d}{dS}(1 - e^{-cS}) = 1, \quad (11.17)$$

or

$$ce^{-cS} = 1. \quad (11.18)$$

Setting $S = 0$ we then find that the transition takes place at $c = 1$.

In other words, the random graph can have a giant component only if $c > 1$. At $c = 1$ and below, we have $S = 0$ and there is no giant component.

This does not entirely solve the problem, however. Technically, we have proved that there can be no giant component for $c \leq 1$, but not that there has to be a giant component for $c > 1$ —in the latter regime there are two solutions for S , one of which is the solution $S = 0$ in which there is no giant component. So which of these solutions is the correct one that describes the true giant component?

In answering this question, we will see another way of thinking about the formation of the giant component. Consider the following process. Let us find a small connected set of nodes somewhere in our network—say a dozen or so—as shown in Fig. 11.3a. In the limit of large n such a set is bound to exist somewhere in the network, so long as the probability of an edge is non-zero. We will divide the set into its *core* and its *periphery*. The periphery is the nodes that have at least one neighbor outside the set—the lighter gray region in the figure. The core is the nodes that only have connections inside the set—the darker gray.

Now imagine enlarging our set by adding to it all those nodes that are immediate neighbors, connected by at least one edge to the set—see Fig. 11.3b. Now the old periphery is part of the core and there is a new periphery consisting of the nodes just added. How big is this new periphery? We know that each node in the old periphery is connected with independent probability p to every other node. If there are s nodes in our set, then there are $n - s$ nodes outside the set, and the average number of connections a node in the periphery has to outside nodes is

$$p(n - s) = c \frac{n - s}{n - 1} \simeq c, \quad (11.19)$$

where the equality becomes exact in the limit $n \rightarrow \infty$. This means that the average number of immediate neighbors of the set—the expected size of the new periphery when we grow the set—is c times the size of the old periphery.³

³It could happen that two nodes in the set have the same neighbor outside the set, in which case our calculation would overcount the number of neighbors. The probability of this happening,

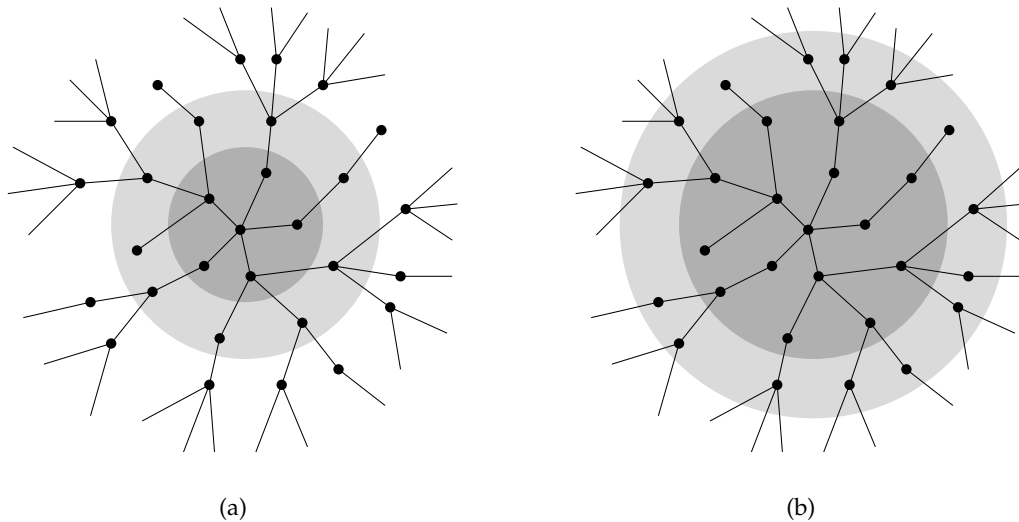


Figure 11.3: Growth of a node set in a random graph. (a) A set of nodes (inside the gray circles) consists of a core (dark gray) and a periphery (lighter). (b) If we grow the set by adding to it those nodes immediately adjacent to the periphery, then the periphery nodes become a part of the new core and a new periphery is added.

We can repeat this argument as many times as we like, growing the set again and again, and each time the average size of the periphery will increase by a factor of c . Thus if $c > 1$ the average size of the periphery will grow exponentially. On the other hand, if $c < 1$ it will shrink exponentially and eventually dwindle to zero.⁴ If it grows exponentially, our connected set of nodes will eventually form a component comparable in size to the whole network—a giant component—while if it dwindles the set will only ever have finite size and no giant component will form.

So we see that indeed we expect a giant component if (and only if) $c > 1$.

however, is very small when n is large: if there are many nodes outside the set, and given that each node chooses its neighbors at random, the chances of two nodes inside the set having the same neighbor outside the set are vanishingly small. This observation is related to the fact that the networks generated by random graph models are “locally tree-like,” as discussed in detail in Section 12.4.

⁴The marginal case where c is exactly equal to 1 is more complicated and we will not study it here. A detailed analysis shows that when $c = 1$ there is technically no giant component, in the sense of a component whose size increases as n , but there is a largest component with size scaling as $n^{2/3}$ [70].

When there is a giant component the size of that component will be given by the larger solution to Eq. (11.16). This now allows us to calculate the size of the giant component for all values of c . The results are shown in Fig. 11.2b. As the figure shows, the size of the giant component grows rapidly from zero as the value of c passes 1, and tends towards $S = 1$ as c becomes large. (We have to solve for the solution of Eq. (11.16) numerically, since the equation has no closed-form solution, but this is easy enough to do.)

11.5.1 CAN THERE BE MORE THAN ONE GIANT COMPONENT?

So far we have been assuming that there is only one giant component in our network, only one component whose size scales with n . This seems plausible, particularly given our experience in Chapter 10 with real-world networks, which usually have only one large component, but might it be possible to have two or more giant components in a network? In Section 10.1 we gave a rough argument suggesting that there can only be one giant component. We are now in a position to make this argument more rigorous for the case of the random graph.

Imagine that we generate a random graph in the usual way by placing edges with probability $p = c/(n-1)$, and suppose as usual that the network is sparse, meaning that c grows slower than n . To be concrete, let us say that c grows no faster than n^a for large n and some positive constant $a < 1$. In the common case where c is constant, for instance, any value of a between zero and one would work.

Now suppose that we also add a further set of edges to the network by going through all node pairs that don't yet have an edge and placing an edge between them with a different probability $p' = c/(n-1)^{1+a}$. The end result is a random graph again, but now with a larger edge probability, equal to $p + p'$ when n is large, and an average degree given by

$$\begin{aligned} c' &= (n-1)(p + p') = (n-1) \left[\frac{c}{n-1} + \frac{c}{(n-1)^{1+a}} \right] \\ &= c \left[1 + \frac{1}{(n-1)^a} \right]. \end{aligned} \quad (11.20)$$

But as $n \rightarrow \infty$ the final term vanishes since a is positive, and we get simply $c' = c$. Hence in this limit we have a random graph with the same mean degree as before. In effect, what we have done is sprinkled a few extra edges over our random graph, but with a density that vanishes in the limit of large n , so that the resulting ensemble of networks is unchanged from the one we had before (in this limit). To put that another way, this two-step process, of first placing

edges with probability p then placing more edges with probability p' , is just a complicated way of generating an ordinary random graph with mean degree c . So if we can prove that networks generated this way have no more than one giant component, then we will have the result we were looking for.

In order for such a network to have two or more giant components the components must be separate at both steps of the generating process: they must be separate after the first round of edges are placed with probability p and they must remain separate after the second round are placed with probability p' . Suppose then that there are two or more separate giant components after the first round and take any two of those giant components, with sizes S_1n and S_2n , where S_1 and S_2 are the fractions of the network filled by each. The number of distinct pairs of nodes i, j such that i is in the first giant component and j is in the second is $S_1n \times S_2n = S_1S_2n^2$ and by definition none of those pairs is connected by an edge, since if they were the two components would be one. Now we sprinkle our extra edges with probability p' and if the two giant components are to remain separate we require that none of the $S_1S_2n^2$ node pairs get connected by an added edge. The probability q of this happening is $q = (1 - p')^{S_1S_2n^2}$. Taking logs we get

$$\begin{aligned} \ln q &= S_1S_2n^2 \ln(1 - c/(n-1)^{1+a}) \\ &= S_1S_2n^2 \left[-\frac{c}{(n-1)^{1+a}} - \frac{c^2}{2(n-1)^{2+2a}} - \dots \right] \\ &\simeq -cS_1S_2n^{1-a}, \end{aligned} \tag{11.21}$$

where the approximate equality becomes exact in the limit of large n . Taking exponentials again, we get

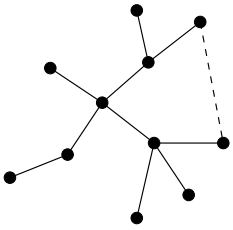
$$q = e^{-cS_1S_2n^{1-a}}, \tag{11.22}$$

which goes to zero as $n \rightarrow \infty$ since $a < 1$. Given, as we have said, that our final random graph is a true random graph with average degree c , this now establishes the result we wanted: in the limit of large n , the probability that we will have two separate giant components in such a network goes to zero.

11.6 SMALL COMPONENTS

We have seen that in a random graph with $c > 1$ there exists a giant component whose size grows in proportion to n and which fills an extensive fraction of the network. That fraction is typically less than 100%, however. What is the structure of the remainder of the network? It cannot contain additional giant components—there is only one giant component, as we have seen—so it must

Recall that a tree is a network or subnetwork that has no loops—see Section 6.8.



If we add an edge (dashed) to a tree we create a loop.

be made up of components whose size grows slower than n . We call these the *small components*. Since their size grows slower than n there must necessarily be many of them as n becomes large (so long as the giant component doesn't fill the whole network) and they can in general have a range of sizes, some larger than others. It is possible to calculate the entire distribution of these sizes, but we will not do so here—the calculation is lengthy and the result is not particularly illuminating.⁵ What is interesting, however, is the average size of a small component, which is quite easy to calculate.

The crucial insight that makes the calculation possible is that the small components are trees, which we can demonstrate by the following argument. Consider a small component of s nodes that takes the form of a tree. As shown in Section 6.8, a tree of s nodes contains $s - 1$ edges, which is the smallest number of edges that such a set of nodes can have and still be connected together. If we add any other edge to the component then we will create a loop (since we will be adding a new path between two nodes that are already connected—see figure) and hence the component will no longer be a tree. In a Poisson random graph the probability of such an edge is the same as for any other edge, $p = c/(n - 1)$. The total number of places where we could add an extra edge to the component is given by the number of distinct pairs of nodes minus the number that are already connected by an edge, or

$$\binom{s}{2} - (s - 1) = \frac{1}{2}(s - 1)(s - 2), \quad (11.23)$$

and so the average total number of extra edges in the component is $\frac{1}{2}(s - 1)(s - 2) \times c/(n - 1)$. For any given value of s this number tends to zero as $n \rightarrow \infty$, and hence there are no loops in the component and the component is a tree.

Given this observation, consider now a node i in a small component of a random graph, as depicted in Fig. 11.4a. Each of i 's edges leads to a separate subnetwork—the shaded regions in the figure—and because the whole component is a tree we know that these subnetworks are not connected to one another, other than via node i , since if they were there would be a loop in the component and it would not be a tree. Thus, the size of the component to which i belongs is the sum of the sizes of the subnetworks reachable along each of its edges, plus 1 for node i itself. Let us call the sizes of the subnetworks $t_1 \dots t_k$, where k is the degree of node i (which is 3 in the figure, but could be anything we like).

⁵However, in Section 12.10.9 we calculate the complete distribution of small component sizes for a more general model, the configuration model, that includes the Poisson random graph as a special case.

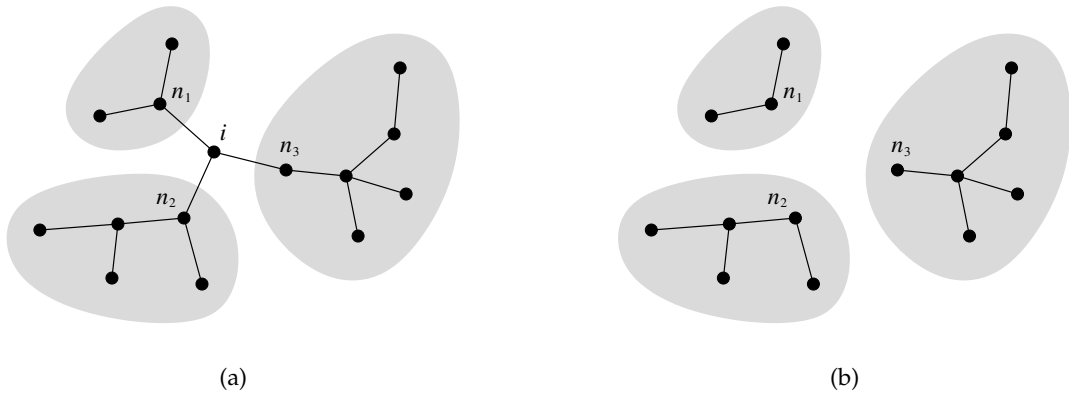


Figure 11.4: The size of a small component in a random graph. (a) The size of the component to which a node i belongs is the sum of the number of nodes in each of the subcomponents (shaded regions) reachable via i 's neighbors n_1, n_2, n_3 , plus one for i itself. (b) If node i is removed the subcomponents become components in their own right.

Then the size of the component is

$$s = 1 + \sum_{m=1}^k t_m. \quad (11.24)$$

To calculate the mean size of a small component we now want to average this expression over many different nodes in the small components of the network. We will perform this average in two stages. First, we average over only the nodes that have degree k . Taking the average of both sides of Eq. (11.24) we get

$$\langle s \rangle_k = 1 + \sum_{m=1}^k \langle t_m \rangle, \quad (11.25)$$

where the subscript k reminds us that we are averaging over nodes of degree k and $\langle t_m \rangle$ is the average size of the subnetwork that the m th neighbor of a node belongs to. But all neighbors are equivalent in a random graph—there is nothing to distinguish the m th neighbor from any other—and hence $\langle t_m \rangle$ has the same value for all m , which we will write as simply $\langle t \rangle$. Thus

$$\langle s \rangle_k = 1 + k \langle t \rangle. \quad (11.26)$$

Now we average this expression further, not just over nodes with degree k but over small-component nodes with any degree, which gives the average

size $\langle s \rangle$ of the component to which such a node belongs thus:

$$\langle s \rangle = 1 + \langle k \rangle_{\text{small}} \langle t \rangle, \quad (11.27)$$

where $\langle k \rangle_{\text{small}}$ is the average degree of a node in a small component.

It remains for us to calculate the values of $\langle k \rangle_{\text{small}}$ and $\langle t \rangle$. The first is straightforward, once we notice that the average degree of a node in a small component is not equal to the average degree c in the network as a whole. A node in a small component can only be connected to other nodes in small components. The giant component (if there is one) fills a fraction S of the network, meaning that the small components fill a fraction $1 - S$ and hence there are $(1 - S)n$ nodes in small components. Each of these nodes is connected to any of the others with the usual probability p , and hence the average degree of a node in a small component is

$$\langle k \rangle_{\text{small}} = [(1 - S)n - 1]p = [(1 - S)n - 1] \frac{c}{n - 1} \simeq (1 - S)c, \quad (11.28)$$

where we have used Eq. (11.6) and the approximate equality becomes exact in the limit of large n . In other words, the average degree of a node in a small component is smaller than the average degree of a node in the network as a whole, by a factor of $1 - S$.

And what about the value of $\langle t \rangle$? To determine this, let us return to Fig. 11.4 and consider a slightly modified network, the network in which node i is removed, along with all its edges, as shown in panel (b) of the figure.⁶ This network is still a random graph with the same value of p —each possible edge is still present with independent probability p —but the number of nodes has decreased by one, from n to $n - 1$. In the limit of large n , however, this decrease is negligible. The average properties, such as size of the giant component and sizes of small components, will be indistinguishable for random graphs with sizes n and $n - 1$ but the same p .

In this modified network, what were previously the subnetworks containing the neighbors n_1, n_2, \dots of node i (the shaded regions) are now separate small components in their own right. But since the network is still a random graph with the same edge probability p as the original network, the average size $\langle t \rangle$ of the component that the nodes n_1, n_2, \dots belong to is simply equal to the average size of *any* small component—there is nothing special about either the components or the nodes that would say otherwise. In other words $\langle t \rangle = \langle s \rangle$.

⁶This trick of removing a node is called a *cavity method*. Cavity methods are used widely in physics for the solution of all kinds of problems and are a powerful method for many calculations on lattices and in low-dimensional spaces as well as on networks [332].

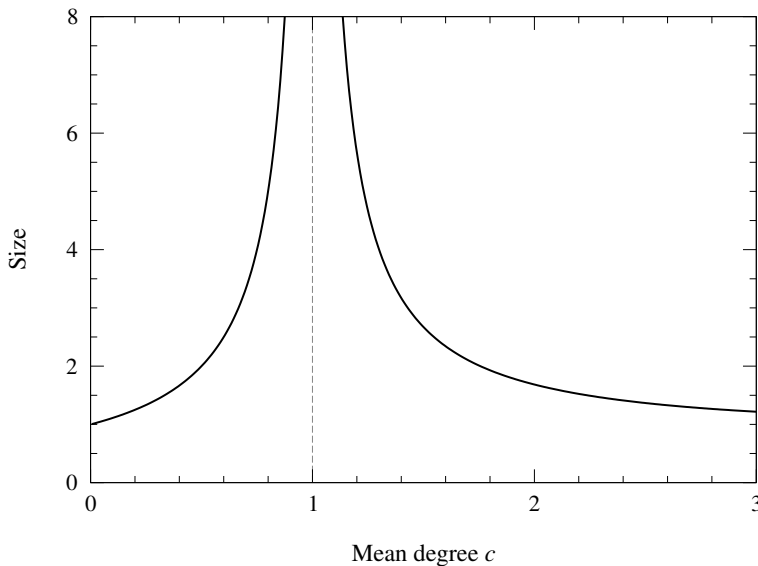


Figure 11.5: Average size of the small components in a random graph. The average size $\langle s \rangle$ of the component to which a randomly chosen node in a small component belongs, calculated from Eq. (11.29).

Putting this result together with Eqs. (11.27) and (11.28), we now have $\langle s \rangle = 1 + (1 - S)c\langle s \rangle$, or

$$\langle s \rangle = \frac{1}{1 - c + cS}. \quad (11.29)$$

This is the average size of the small component to which a randomly chosen node in a random graph belongs.⁷

When $c < 1$ and there is no giant component, Eq. (11.29) gives simply $\langle s \rangle = 1/(1 - c)$. When there is a giant component evaluating $\langle s \rangle$ is more complicated, because we must first solve for S before finding the value of $\langle s \rangle$, but the calculation can still be done—we solve Eq. (11.16) for S and then substitute into Eq. (11.29). Figure 11.5 shows a plot of the value of $\langle s \rangle$ as a function of c .

Note how $\langle s \rangle$ diverges when $c = 1$. (At this point $S = 0$, so the denominator

⁷Note that this is not the same thing as the average size of a small component. Since there are more nodes in a component of larger size than in a component of smaller size, the mean size of the component to which a node belongs is a biased average, giving more weight to larger components. For most practical purposes, however, Eq. (11.29) turns out to be the most useful metric of average component size and it is the one we use in this book.

of (11.29) vanishes.) Thus, if we slowly increase the mean degree c of our network from some small initial value less than 1, the average size of the component to which a node belongs gets bigger and bigger and finally becomes infinite exactly at the point $c = 1$ where the giant component appears. For $c > 1$ Eq. (11.29)—which measures only the sizes of the non-giant components—tells us that these components get smaller again as c increases. Thus the general picture we have is one in which the small components get larger up to $c = 1$, where they diverge and the giant component appears, then smaller again as the giant component grows larger.

Although the random graph is certainly not a realistic model of most networks, this general picture of the component structure of the network turns out to be a good guide to the behavior of networks in the real world. If a network has a low density of edges, then typically it consists only of small components, but if the density becomes high enough then a single large component forms, usually accompanied by some separate small components. This is a good example of the way in which simple models of networks can give us a feel for how more complicated real-world systems should behave.

11.7 PATH LENGTHS

In Sections 4.6 and 10.2 we discussed the small-world effect, the observation that the typical lengths of paths between nodes in networks tend to be short. We can use the random graph model to shed light on how the small-world effect arises by examining the behavior of the network diameter in the model.

Recall that the diameter of a network is the longest distance between any two nodes in the same component—the “longest shortest path,” if you like. As we now show, the diameter of a random graph varies with the number n of nodes as $\ln n$. Since $\ln n$ is a relatively small number even when n is large, this offers some explanation of the small-world effect, although, as we will see, it also leaves some questions open.

The basic idea behind the calculation of the diameter of a random graph is straightforward. In Section 11.5 we argued that if we grow a set of nodes in a random graph by repeatedly adding the neighbors of the set to it, the number of neighbors added goes up by a factor of c , on average, on each step. Imagine growing such a set starting from a single node and working outward. The average number of nodes one step away is clearly c , the average degree, and if it grows by a factor of c on each additional step then the average number of nodes s steps away must be c^s . Since this expression grows exponentially with s it doesn't take very many such steps before the number of nodes reached is equal to the total number of nodes in the whole network—this happens when

See Section 6.11.1 for a discuss of shortest distances and diameters.

$c^s \simeq n$ or equivalently $s \simeq \ln n / \ln c$. At this point, roughly speaking, every node is within s steps of our starting point, implying that the diameter of the network is approximately $\ln n / \ln c$.

Although the random graph is, as we have said, not an accurate model of most real-world networks, this is, nonetheless, believed to be the basic mechanism behind the small-world effect: the number of nodes within distance s of a particular starting point grows exponentially with s and hence the diameter is logarithmic in n . We discuss the comparison with real-world networks in more detail in Section 11.8.

The argument above is only approximate. First, it really calculates the “radius” of the network, not the diameter—it tells us the maximum distance from an average starting point to other nodes, not the maximum distance in the network as a whole. Moreover, while it’s true that there are on average c^s nodes s steps away from any starting point so long as s is small, this result must break down once c^s becomes comparable with n , since clearly the number of nodes at distance s cannot exceed the number of nodes in the whole network. (Indeed it cannot exceed the number in the giant component.)

We can deal with both of these problems by considering two different starting nodes i and j . The average numbers of nodes s and t steps from these starting nodes will then be c^s and c^t so long as we stay in the regime where both these numbers are much less than n . In the following calculation we consider only configurations in which both remain smaller than order n in the limit $n \rightarrow \infty$ so as to satisfy this condition.

The situation we consider is depicted in Fig. 11.6, with the two nodes i and j each surrounded by a “ball” or neighborhood consisting of all nodes with distances up to and including s and t , respectively. If there is an edge between the “surface” (i.e., most distant nodes) of one neighborhood and the surface of the other, as depicted by the dashed line, then the length d_{ij} of the shortest path from i to j is at most $s + t + 1$. Conversely, if there is no such edge, if the two balls have not met yet, then d_{ij} must be greater than $s + t + 1$. Or, to put that another way, the probability $P(d_{ij} > s + t + 1)$ that d_{ij} is greater than $s + t + 1$ is equal to the probability that there is no edge between the two surfaces.

There are on average $c^s \times c^t$ pairs of nodes such that one lies on each surface, and each pair is connected with probability $p = c/(n - 1) \simeq c/n$ (assuming n to be large) or not with probability $1 - p$. Hence $P(d_{ij} > s + t + 1) = (1 - p)^{c^{s+t}}$. Defining for convenience $\ell = s + t + 1$, we can also write this as

$$P(d_{ij} > \ell) = (1 - p)^{c^{\ell-1}} = \left(1 - \frac{c}{n}\right)^{c^{\ell-1}}. \quad (11.30)$$

This is the same trick employed by the two-source breadth-first search algorithm of Section 8.5.4 to avoid probing all the nodes in a network—growing outwards from two points until you meet in the middle.

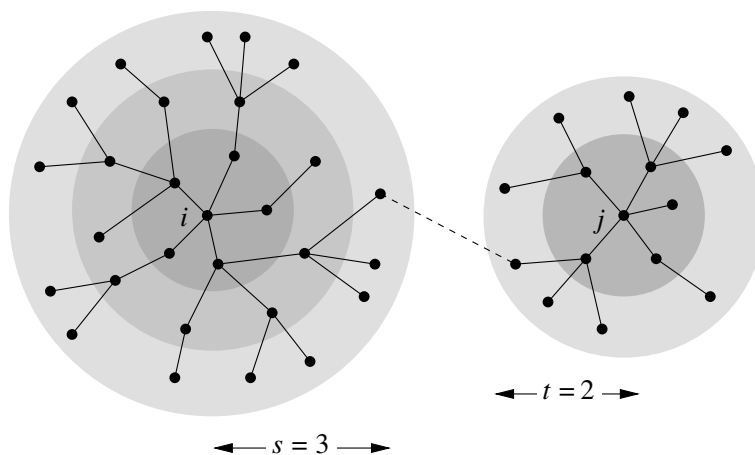


Figure 11.6: Neighborhoods of two nodes in a random graph. In the argument given in the text we consider the sets of nodes within distances s and t respectively of two randomly chosen nodes i and j . If there is an edge between any node on the surface of one neighborhood and any node on the surface of the other (dashed line), then there is a path between i and j of length $s + t + 1$.

Taking logs of both sides, we find that

$$\ln P(d_{ij} > \ell) = c^{\ell-1} \ln\left(1 - \frac{c}{n}\right) \approx -\frac{c^\ell}{n}, \quad (11.31)$$

where the approximate inequality becomes exact as $n \rightarrow \infty$. Thus in this limit

$$P(d_{ij} > \ell) = \exp\left(-\frac{c^\ell}{n}\right). \quad (11.32)$$

The diameter of the network is the smallest value of ℓ such that $P(d_{ij} > \ell)$ is zero, i.e., the value such that no matter which pair of nodes we happen to pick there is zero chance that they will be separated by a distance greater than ℓ . In the limit of large n , Eq. (11.32) will tend to zero only if c^ℓ grows faster than n , meaning that our smallest value of ℓ is the value such that $c^\ell = an^{1+\epsilon}$ with a constant and $\epsilon \rightarrow 0$ from above. Note that we can, as promised, achieve this while keeping both c^s and c^t smaller than order n , so that our argument remains valid. For instance, since $c^{s+t} = c^{\ell-1} = (a/c)n^{1+\epsilon}$, we could choose both c^s and c^t growing as $n^{(1+\epsilon)/2}$ on average.

Rearranging $c^\ell = an^{1+\epsilon}$ for ℓ , we now find our expression for the diameter:

$$\ell = \frac{\ln a}{\ln c} + \lim_{\epsilon \rightarrow 0} \frac{(1+\epsilon) \ln n}{\ln c} = A + \frac{\ln n}{\ln c}, \quad (11.33)$$

where A is a constant.⁸ Apart from the constant, this is the same result as we found previously using a rougher argument. The constant is known—it has a rather complicated value in terms of the Lambert W -function [172]—but for our purposes the important point is that it is (asymptotically) independent of n . Thus, the diameter indeed increases only slowly with n , as $\ln n$, making it relatively small in large random graphs.

The logarithmic dependence of the diameter on n offers some explanation of the small-world effect discussed in Section 4.6. Even in a network such as the acquaintance network of the entire world, with over seven billion inhabitants (at the time of writing), the value of $\ln n / \ln c$ can be quite small. Supposing each person to have about a thousand acquaintances, we would get

$$\ell = \frac{\ln n}{\ln c} = \frac{\ln(7 \times 10^9)}{\ln 1000} = 3.28 \dots, \quad (11.34)$$

which is easily small enough to account for the results of the small-world experiments performed by Milgram and others [142, 333, 447].

In practice Eq. (11.33) appears to be a good guide to the behavior of many real networks. For example, Fig. 11.7 shows average shortest-path distances in the Facebook friendship networks of students at 100 different US universities, from a study by Jacobs *et al.* [246]. The distances are plotted as a function of $\ln n$ and, as the figure shows, approximately lie on a straight line, as implied by Eq. (11.33).

The agreement is not perfect though. There is some scatter among the points, which could be due to variation in the mean degree c between universities, or just to statistical fluctuations. Note also that the figure shows average distance, not diameter (i.e., largest distance) which is what we calculated in Eq. (11.33). The latter provides an upper bound on the former, so we expect the average distance to grow no faster than the diameter—i.e., no faster than $\ln n$ —but we should not expect the two to be equal.

More importantly, however, the random graph is not a very good model of most real-world networks. There are clearly many things wrong with it, as we now discuss.

⁸There are still some holes in our argument. In particular, we have assumed that the product of the numbers of nodes on the surface of our two neighborhoods is c^{s+t} when in practice this is only the average value and there will in general be some variation. We need to prove that this variation is small, i.e., that the result is sufficiently “concentrated” around the average value. Also the calculation should really be confined to the giant component, since the longest path always falls in the giant component in the limit of large n . For a careful treatment of these issues see, for instance, Fernholz and Ramachandran [172].

A thousand appears to be a reasonable figure for the typical person’s number of acquaintances. Bernard and collaborators [54, 55, 261] estimated the average number of acquaintances for people in several cities and found figures ranging from a few hundred to about 2000—see Section 4.2.1.

11.8 PROBLEMS WITH THE RANDOM GRAPH

The Poisson random graph is one of the best studied models of networks. In the decades since its first proposal it has given us a tremendous amount of insight into the expected structure of networks of all kinds, particularly with respect to component sizes and network diameters. The fact that it is both simple to describe and straightforward to study using analytic methods makes it an excellent tool for investigating all sorts of network phenomena. We will return to the random graph many times in the remainder of this book to help us understand the way networks behave.

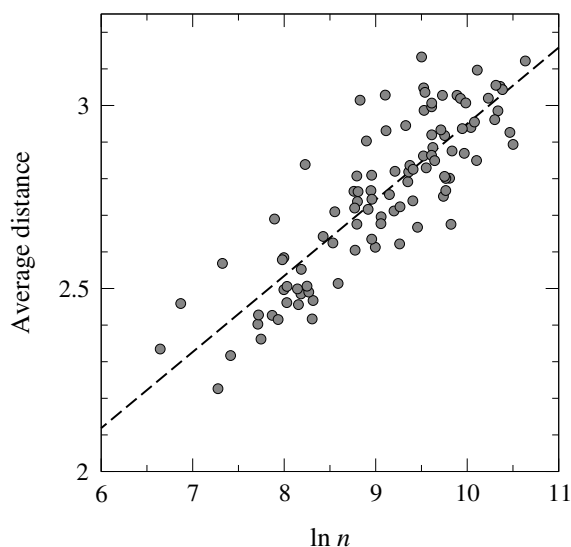


Figure 11.7: Average shortest path distance in Facebook friendship networks. The 100 points in this plot represent average distances in the Facebook friendship networks of students at 100 different US universities, plotted as a function of $\ln n$, the log of the number of nodes in the networks. The dashed line is the best straight-line fit. After Jacobs *et al.* [246].

The random graph does, however, have some severe shortcomings as a network model. There are ways in which it is completely unlike the real-world networks we have seen in previous chapters. One clear problem is that it shows essentially no transitivity or clustering. In Section 11.4 we saw that the clustering coefficient of a random graph is $C = c/(n-1)$, which tends to zero in the limit of large n . And even for the finite values of n appropriate to real-world networks the value of C in the random graph is often very small. For the acquaintance network of the world's population, with its $n \approx 7$ billion people, each having about $c = 1000$ acquaintances, a random graph with the same n and c would have a clustering coefficient of

$$C \approx \frac{1000}{7000000000} \approx 10^{-7}. \quad (11.35)$$

Whether the clustering coefficient of the real acquaintance network is 0.01 or 0.5 hardly matters. (It is probably somewhere in between.) Either way it is clear that the random graph and the true network are in strong disagreement.⁹

The random graph also differs from real-world networks in other ways. For instance, there is no correlation between the degrees of adjacent nodes—necessarily so, since the edges are placed completely at random. The degrees in real networks, by contrast, are usually correlated, as discussed in Section 10.7.

⁹This disagreement between random graphs and real networks, highlighted particularly by Watts and Strogatz [466], was one of the things that prompted the current wave of interest in networks in the mathematical sciences starting in the 1990s.

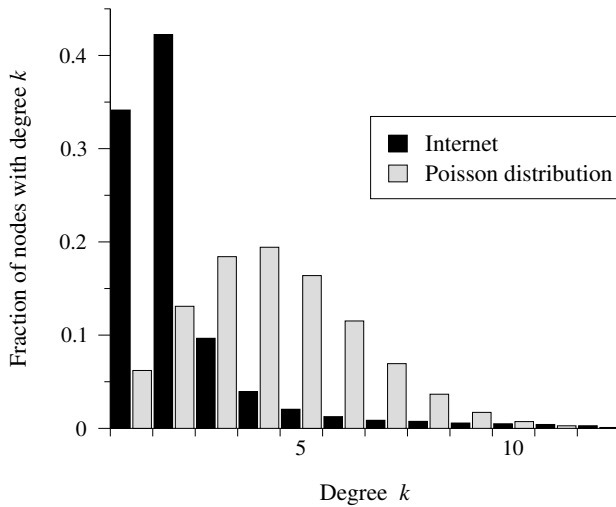


Figure 11.8: Degree distribution of the Internet and a Poisson random graph. The darker bars in this plot show the fraction of nodes with given degrees in a network representation of the Internet at the level of autonomous systems. The lighter bars show the same measure for a random graph with the same average degree as the Internet. Even though the two distributions have the same averages, it is clear that they are entirely different in shape.

Many, perhaps most, real-world networks also show grouping of their nodes into “communities,” as discussed in Chapter 14, but random graphs have no such structure. And there are many other examples of interesting structure in real networks that is absent from the random graph.

However, perhaps the most significant respect in which the properties of random graphs depart from those of real-world networks is the shape of their degree distribution. As discussed in Section 10.3, real networks typically have right-skewed degree distributions, with most nodes having low degree but with a small number of high-degree “hubs” in the tail of the distribution. The random graph, on the other hand, has a Poisson degree distribution, Eq. (11.10), which is not right-skewed to any significant extent. Consider Fig. 11.8, for example, which shows a histogram of the degree distribution of the Internet (darker bars), measured at the level of autonomous systems. The right-skewed form is clearly visible in this distribution. On the same figure we show the Poisson degree distribution of a random graph (lighter bars) with the same average degree c as the Internet example. Despite having the same averages, the two distributions are clearly entirely different. It turns out that this difference

has a profound effect on all sorts of properties of the network—we will see many examples in this book. This makes the Poisson random graph inadequate to explain many of the interesting phenomena we see in networks today, including resilience phenomena, epidemic spreading processes, percolation, and many others.

Luckily it turns out to be possible to generalize the random graph model to allow for non-Poisson degree distributions. This development, which leads to some of the most beautiful results in the mathematics of networks, is described in the next chapter.

EXERCISES

11.1 Consider the random graph $G(n, p)$ with mean degree c .

- Show that in the limit of large n the expected number of triangles in the network is $\frac{1}{6}c^3$. This means that the number of triangles is constant, neither growing nor vanishing in the limit of large n .
- Show that the expected number of connected triples in the network (as defined on page 184) is $\frac{1}{2}nc^2$.
- Hence calculate the clustering coefficient C , as defined in Eq. (7.28), and confirm that it agrees for large n with the value given in Eq. (11.11).

11.2 Consider the random graph $G(n, p)$ with mean degree c .

- Argue that the probability that a node of degree k belongs to a small component is $(1 - S)^k$, where S is the fraction of the network occupied by the giant component.
- Thus, using Bayes' theorem (or otherwise) show that the fraction of nodes in small components that have degree k is $e^{-c}c^k(1 - S)^{k-1}/k!$.

11.3 Write a computer program in the language of your choice that generates a random graph drawn from the model $G(n, m)$ for given values of n and the average degree $c = 2m/n$, then calculates the size of its largest component. Use your program to find the size of the largest component in a random graph with $n = 1\,000\,000$ and $c = 2 \ln 2 = 1.3863\dots$ and compare your answer to the analytic prediction for the giant component of $G(n, p)$ with the same size and average degree. You should find good agreement, even though the models are not identical.

11.4 Consider the random graph $G(n, p)$ with n large.

- If the network has a giant component that fills exactly half of the network, what is the average degree of a node?

- b) For this same random graph what is the probability that a node has degree exactly 5?
- c) What is the probability that a node belongs to the giant component if it has degree exactly 5?
- d) Hence or otherwise, calculate the fraction of nodes in the giant component that have degree exactly 5.

11.5 As we have seen, if a node in a random graph has (at least) one edge connecting it to the giant component then it is itself in the giant component. On the other hand, if a node has at least *two* connections to the giant component then it is in the giant *bicomponent*. (See Section 7.2.3 for a discussion of bicomponents.)

- a) Consider a random graph $G(n, p)$ with mean degree c , and let S be the fraction of the network filled by the giant component. Show that the probability of a node having no connections at all to the giant component is e^{-cS} in the limit of large n .
- b) Show that the probability of exactly one connection to the giant component—no more and no less—is cSe^{-cS} in the limit of large n .
- c) Hence show that the fraction T of the network filled by the giant bicomponent is $T = 1 - (1 + cS)e^{-cS}$.
- d) Show that this can be rewritten as $T = S + (1 - S)\ln(1 - S)$. Hence argue that the giant bicomponent is always smaller than the giant component, unless the giant component fills the whole network or there is no giant component at all.

11.6 Equation (11.28) tells us that the average degree in the small components of a random graph is $(1 - S)c$.

- a) Give an argument to show that the small components on their own, without the giant component (if there is one), themselves constitute a random graph, and hence that the average degree in the small components must be less than 1.
- b) Use this result to construct an alternative argument that there must be a giant component in a random graph if $c > 1$.

11.7 We can make a directed equivalent of the random graph by taking n nodes and placing directed edges with probability p between every pair of distinct nodes. We explicitly and separately consider placing an edge in each direction between every node pair, so a given pair could end up connected by zero edges, one edge, or two edges in opposite directions.

- a) What is the average number m of directed edges in the network in terms of n and p ? Hence what is the average degree c (either in or out) of a node?
- b) If we were to discard the directions of the edges, making an undirected network, what would the average degree be? Hence show that the fraction W of the network occupied by the giant weakly connected component is a solution of $W = 1 - e^{-2cW}$ in the limit of large n .
- c) Let S be the fraction of the network occupied by the giant strongly connected component. In order to belong to the giant strongly connected component, a node must have at least one outgoing edge leading to another node in the giant strongly

connected component and at least one incoming edge leading from a node in the giant strongly connected component. Hence derive an equation analogous to that for W above that must be satisfied by S in the limit of large n .

11.8 The *cascade model* is a simple mathematical model of a directed acyclic network, sometimes used to model food webs. We take n nodes labeled $i = 1 \dots n$ and place an undirected edge between each distinct pair with independent probability p , just as in the ordinary random graph. Then we add directions to the edges such that each edge runs from the node with numerically higher label to the node with lower label. This ensures that all directed paths in the network run from higher to lower labels and hence that the network is acyclic, as discussed in Section 6.4.1.

- Show that the average in-degree of node i in the ensemble of the cascade model is $\langle k_i^{\text{in}} \rangle = (n - i)p$ and the average out-degree is $\langle k_i^{\text{out}} \rangle = (i - 1)p$.
- Show that the expected number of edges that connect to nodes i and lower from nodes above i is $(ni - i^2)p$.
- Assuming n is even, what are the largest and smallest values of this quantity and where do they occur?

In a food web this expected number of edges from high- to low-numbered nodes is a rough measure of energy flow and the cascade model predicts that energy flow will be largest in the middle portions of a food web and smallest at the top and bottom.

11.9 We can make a simple random graph model of a network with clustering or transitivity as follows. We take n nodes and go through each distinct trio of three nodes, of which there are $\binom{n}{3}$, and with independent probability p we connect the members of the trio together using three edges to form a triangle, where $p = c / \binom{n-1}{2}$ with c a constant.

- Show that the mean degree of a node in this model network is $2c$.
- Show that the degree distribution is

$$p_k = \begin{cases} e^{-c} c^{k/2} / (\frac{1}{2}k)! & \text{if } k \text{ is even,} \\ 0 & \text{if } k \text{ is odd.} \end{cases}$$

- Show that the clustering coefficient, Eq. (7.28), is $C = 1/(2c + 1)$.
- Show that when there is a giant component in the network, its expected size S as a fraction of network size satisfies $S = 1 - e^{-cS(2-S)}$.
- What is the value of the clustering coefficient when the giant component fills half the network?

Random graph models with clustering are discussed in Section 12.11.5.

CHAPTER 12

THE CONFIGURATION MODEL

An introduction to the configuration model and related models, which mimic networks with arbitrary degree distributions and other features

IN THE previous chapter we looked at the classic random graph model in which pairs of nodes are connected at random with uniform probability. Although this model has proved tremendously useful as a source of insight into the structure of networks, it also has, as described in Section 11.8, a number of serious shortcomings. Chief among these is its degree distribution, which follows the Poisson distribution and is quite different from the right-skewed degree distributions seen in most real-world networks. In this chapter we introduce the configuration model, a more sophisticated kind of random graph which can have any degree distribution and yet is still exactly solvable for many properties in the limit of large network size.

The configuration model is one of the most important theoretical models in the study of networks. For many purposes it strikes an ideal balance between realism and simplicity, and is frequently the first model that network scientists turn to when studying a new question or process. If you are interested in something that happens on networks—the flow of traffic, the spread of a disease, the evolution of a dynamical system—it is often a good first step to see how it behaves on networks generated using the configuration model.

Also in this chapter we describe briefly a number of additional random graph models that incorporate other features seen in real-world networks, including models of directed and bipartite networks, models of transitivity and assortative mixing, and models for networks that evolve over time.

12.1 THE CONFIGURATION MODEL

We can turn the random graph of Chapter 11 into a much more flexible and powerful model of networks by modifying it so that the degrees of its nodes are no longer restricted to having a Poisson distribution. In fact, it turns out to be possible to modify the model so as to give the network any degree distribution we please. Just as with the Poisson random graph, which can be defined in several slightly different ways, there is more than one way to define random graphs with general degree distributions. Here we describe two of them, which are roughly the equivalent of the $G(n, m)$ and $G(n, p)$ random graphs of Section 11.1.

The most widely studied of the generalized random graph models is the *configuration model*.¹ The configuration model is actually a model of a random graph with a given degree *sequence*, rather than degree distribution. That is, the exact degree of each individual node in the network is specified, rather than merely the probability distribution from which those degrees are chosen. This in turn fixes the number of edges in the network, since the number of edges is given by Eq. (6.13) to be $m = \frac{1}{2} \sum_i k_i$. Thus this model is in some ways analogous to the $G(n, m)$ random graph model, which also fixes the number of edges in the network. It is quite simple, however, to modify the configuration model for cases where only the degree distribution is known and not the exact degree sequence. We describe how this is done at the end of this section.

Suppose then that we specify the degree k_i of each node $i = 1 \dots n$ in our network. We can create a random network with these degrees using the process depicted in Fig. 12.1. We give each node i a total of k_i “stubs” of edges, also sometimes called “half-edges.” There are $\sum_i k_i = 2m$ stubs in total, where m is the total number of edges. Now we choose two of the stubs uniformly at random and connect them together to form an edge, as indicated by the dashed line in the figure. Then we choose another pair from the remaining $2m - 2$ stubs, connect those, and so on until all the stubs are used up. The end result is a network in which every node has exactly the desired degree.

More specifically, the end result is a particular *matching* of the stubs, a particular set of pairings of stubs with other stubs. The process above generates each possible matching of stubs with equal probability. Technically, the configuration model is defined as an en-

See Section 10.3 for a discussion of the distinction between degree sequences and degree distributions.

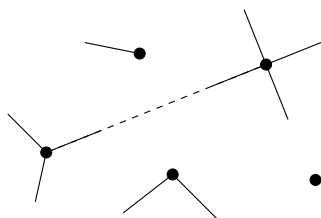


Figure 12.1: The configuration model. Each node is given a number of “stubs” of edges equal to its desired degree. Then pairs of stubs are chosen at random and connected together to form edges (dotted line).

¹The name has its origins in the work of mathematician Béla Bollobás, who, in one of the earliest papers on the topic [69], used the term “configuration” to refer to arrangements of edges in the model.

semble of matchings in which each matching with the chosen degree sequence appears with the same probability and those with any other degree sequence have probability zero. The process above is then a process for drawing networks from the configuration model ensemble.

The uniform distribution over matchings in the configuration model has the important consequence that any stub in a configuration model network is equally likely to be connected to any other. This, as we will see, is the crucial property that makes the model solvable for many of its properties.

There are a couple of catches with the network generation process described here. First, there must be an even number of stubs overall if we want to end up with a network consisting only of nodes and edges, with no dangling stubs left over. This means that the sum $\sum_i k_i$ of the degrees must add up to an even number. We will assume that the degrees we have chosen satisfy this condition; otherwise it is not possible to create a network with the given degree sequence.

A second issue is that the network might contain self-edges or multiedges, or both. There is nothing in the network generation process that prevents us from creating an edge that connects a node to itself or that connects two nodes that are already connected by another edge. One might imagine that one could avoid this by rejecting the creation of any such edges during the matching process, but it turns out that this is not a good idea. A network so generated is no longer drawn uniformly from the set of possible matchings, which means that properties of the model can no longer be calculated analytically, at least by any means currently known. It can also mean that the network creation process breaks down completely. Suppose, for example, that we come to the end of the process, when there are just two stubs left to be joined, and find that those two both belong to the same node so that joining them would create a self-edge. Then either we create the self-edge or the network generation process fails.

In practice, therefore, it makes more sense to allow the creation of both multiedges and self-edges in our networks and the standard configuration model does so. Although some real-world networks have multiedges or self-edges in them, most do not, and to some extent this makes the configuration model less satisfactory as a network model. However, as shown in Section 12.1.1, the density of self-edges and multiedges in the configuration model tends to zero as the network becomes large, which means that their effect on our calculations will typically be negligible so long as we work with reasonably large networks.

As mentioned earlier, we are sometimes (indeed often) interested in the case where it is the degree distribution of the network that is specified rather than the degree sequence. That is, we specify the probability distribution p_k from which the degree sequence is drawn rather than the sequence itself. We

See Section 6.1 for a definition and discussion of self-edges and multiedges.

can define an obvious extension of the configuration model to this case: we draw a degree sequence from the specified distribution and then generate a network with that degree sequence using the technique described above. More precisely, we define an ensemble in which each degree sequence $\{k_i\}$ appears with probability $\prod_i p_{k_i}$ and the probability of a particular matching of stubs is equal to the probability of generating the corresponding degree sequence times the probability of the matching within the standard configuration model.

One small catch with this model is that we must make sure that the degrees k_i that we generate add to an even number. Otherwise, as discussed earlier, we cannot match all the stubs to make the final network—there will always be one left over. This is not a big problem however. If we find that we have generated degrees that add to an odd number, we just throw them away and generate another set.

In practice the difference between the two models is not actually very great. As we will see, the crucial parameter that enters into most of our configuration model calculations is the fraction of nodes that have each possible degree k . In the extended model above, this fraction is, by definition, equal to p_k in the limit of large n . If, on the other hand, the degree sequence is fixed, then we simply calculate the fraction from the degree sequence. In either the case the formulas for calculated quantities are the same (in the limit of large n).

Two important special cases of the model with specified degree distribution are the cases with Poisson and power-law distributions. If we choose a Poisson distribution of node degrees and then generate the corresponding configuration model network, we recover—very nearly—the standard random graph $G(n, p)$. The two are not quite the same since the configuration model can, as we have said, contain multiedges and self-edges, while $G(n, p)$, as normally defined, cannot. Since the density of multiedges and self-edges is small, however, the difference can often be ignored and, in particular, most properties of the configuration model and the Poisson random graph are the same in the limit of large n . Although Poisson degree distributions are rare in real-world networks, the Poisson case can serve as a useful check on our calculations when working with the configuration model. If we derive a result and want to check that it is correct, we can look at the special case for a Poisson degree distribution and see whether we recover the correct answer for $G(n, p)$.

A power-law degree distribution provides a more interesting special case, and one that we will return to repeatedly throughout this book. As discussed in Section 10.4, many networks are observed to have power-law degree distributions, a property that gives rise to some surprising effects. As we will see, we can shed light on these effects by studying the case of the configuration model with power-law degrees.

12.1.1 EDGE PROBABILITY IN THE CONFIGURATION MODEL

A central property of the configuration model is the probability p_{ij} of the occurrence of an edge between two specified nodes, i and j . If either node i or node j has degree zero, then the probability of an edge is also zero, so let us assume that the degrees k_i, k_j are non-zero. Now consider any one of the stubs that emerges from node i . What is the probability that this stub is connected by an edge to any of the stubs of node j ? There are $2m$ stubs in total in the network, or $2m - 1$ excluding the one connected to i that we are currently looking at. Of those $2m - 1$, exactly k_j of them are attached to node j . So, given that any stub in the network is equally likely to be connected to any other, the probability that our particular stub is connected to one of those around node j is $k_j/(2m - 1)$. But there are k_i stubs around node i , so the total probability of a connection between i and j is

$$p_{ij} = \frac{k_i k_j}{2m - 1}. \quad (12.1)$$

Technically, since we have added together the probabilities for all the stubs at node i , this is the expected number of edges between i and j rather than the total edge probability. But in the limit of large m , the value becomes small (for given k_i, k_j) and the expected number of edges and the probability of an edge become equal. Also in the limit of large m we can ignore the -1 in the denominator and hence we can write

$$p_{ij} = \frac{k_i k_j}{2m}. \quad (12.2)$$

This is the form in which this probability is most commonly written. Note that, even though we have assumed $k_i, k_j > 0$, Eq. (12.2) also gives the correct result if either degree is zero, namely that the probability of connection is also zero.

We can use Eq. (12.2) to calculate, for example, the probability of having two edges between the same pair of nodes. The probability of having one edge between nodes i and j is p_{ij} as above. Once we have one edge between the nodes the number of available stubs at each is reduced by one, and hence the probability of having a second edge is given by Eq. (12.2) but with k_i and k_j each reduced by one: $(k_i - 1)(k_j - 1)/2m$. Thus the probability of having (at least) two edges, i.e., of having a multiedge between i and j , is $k_i k_j (k_i - 1)(k_j - 1)/(2m)^2$. Summing this probability over all nodes and dividing by two (to avoid double counting of node pairs), we find that the expected total number of multiedges

in the network is

$$\begin{aligned} \frac{1}{2(2m)^2} \sum_{ij} k_i k_j (k_i - 1)(k_j - 1) &= \frac{1}{2\langle k \rangle^2 n^2} \sum_i k_i (k_i - 1) \sum_j k_j (k_j - 1) \\ &= \frac{1}{2} \left[\frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} \right]^2, \end{aligned} \quad (12.3)$$

where

$$\langle k \rangle = \frac{1}{n} \sum_i k_i, \quad \langle k^2 \rangle = \frac{1}{n} \sum_i k_i^2 \quad (12.4)$$

are the first and second moments of the degree distribution and we have used $2m = \langle k \rangle n$ (see Eq. (6.15)). Thus the expected number of multiedges remains constant as the network grows larger, so long as $\langle k^2 \rangle$ is constant and finite, and the density of multiedges, meaning the number per node, vanishes as $1/n$. We used this result earlier to argue that multiedges in the configuration model are normally rare enough to be ignored.²

Another way to derive the expression in Eq. (12.1) is to observe that there are $k_i k_j$ possible ways to pick a stub from node i and a stub from node j to form an edge. The total number of ways to pick a pair of stubs from the $2m$ stubs in the entire network is $\binom{2m}{2} = m(2m - 1)$. The probability that upon picking a pair of stubs at random we happen to create an edge between i and j is then given by the ratio of these two numbers, which is $k_i k_j / m(2m - 1)$, and when we pick a total of m pairs in succession to create the whole network the expected number of edges we create between i and j is m times this ratio, which gives us Eq. (12.1) again.

The only case for which this derivation is not quite right is the case of self-edges. For a self-edge the number of pairs of stubs is not $k_i k_j$ but instead is $\binom{k_i}{2} = \frac{1}{2} k_i (k_i - 1)$ and hence the probability of a self-edge from node i to itself is $k_i (k_i - 1) / 2(2m - 1)$, or

$$p_{ii} = \frac{k_i (k_i - 1)}{4m}, \quad (12.5)$$

when m is large enough that $2m - 1$ can be safely approximated by just $2m$. We can use this result to calculate the expected number of self-edges in the network, which is given by the sum

$$\sum_i p_{ii} = \sum_i \frac{k_i (k_i - 1)}{4m} = \frac{\langle k^2 \rangle - \langle k \rangle}{2\langle k \rangle}, \quad (12.6)$$

²For networks with power-law degree distributions $\langle k^2 \rangle$ diverges, as described in Section 10.4.2, and in that case the density of multiedges may not vanish or may do so more slowly than $1/n$.

This expression remains constant as $n \rightarrow \infty$ provided $\langle k^2 \rangle$ remains constant, and hence, as with the multiedges, the density of self-edges in the network vanishes as $1/n$ in the limit of large network size.

We can use Eqs. (12.2) and (12.5) to calculate a number of other properties of nodes in the configuration model. For instance, we can calculate the expected number n_{ij} of common neighbors that nodes i and j share. The probability that i is connected to another node l is p_{il} and the probability that j is connected to the same node would likewise normally be p_{jl} . However, as with the calculation of multiedges above, if we already know that i is connected to l , then the number of available stubs at node l is reduced by one and, rather than being given by the normal expression (12.2), the probability of a connection between j and l is then $k_j(k_l - 1)/2m$. Multiplying the probabilities for the two edges and summing over l we then get our expression for the expected number of common neighbors of i and j :

$$\begin{aligned} n_{ij} &= \sum_l \frac{k_i k_l}{2m} \frac{k_j(k_l - 1)}{2m} = \frac{k_i k_j}{2m} \frac{\sum_l k_l(k_l - 1)}{n \langle k \rangle} \\ &= p_{ij} \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle}. \end{aligned} \quad (12.7)$$

Thus, the probability of sharing a common neighbor is equal to the probability $p_{ij} = k_i k_j / 2m$ of having a direct connection times a multiplicative factor that depends only on the mean and variance of the degree distribution but not on the properties of the nodes i and j themselves.³

12.1.2 RANDOM GRAPHS WITH GIVEN EXPECTED DEGREE

The configuration model of the previous section is, as we have said, similar in some ways to the standard random graph $G(n, m)$ described in Section 11.1, in which we distribute a fixed number m of edges at random between n nodes. In the configuration model the total number of edges is again fixed, having value $m = \frac{1}{2} \sum_i k_i$, but in addition we now also fix the individual degree of every node as well.

It is natural to ask whether there is also an equivalent of $G(n, p)$ —the model in which only the probability of edges is fixed and not their number—and indeed there is. We simply place an edge between each pair of nodes i, j

³In this calculation we have ignored the fact that the probability of self-edges, Eq. (12.5), is different from the probability for other edges. As we have seen, however, the density of self-edges in the configuration model tends to zero as $n \rightarrow \infty$, so in that limit it is usually safe to make the approximation that Eq. (12.2) applies for all i and j .

with independent probabilities taking the form of Eq. (12.2). We define a parameter c_i for each node and then place an edge between nodes i and j with probability $p_{ij} = c_i c_j / 2m$. As with the configuration model, we must allow self-edges if the model is to be tractable, and again self-edges must be treated a little differently from ordinary edges. It turns out that the most satisfactory definition of the edge probability is⁴

$$p_{ij} = \begin{cases} c_i c_j / 2m & \text{for } i \neq j, \\ c_i^2 / 4m & \text{for } i = j, \end{cases} \quad (12.8)$$

where m is now defined by⁵

$$\sum_i c_i = 2m. \quad (12.9)$$

With this choice the average number of edges in the network is

$$\sum_{i < j} p_{ij} + \sum_i p_{ii} = \sum_{i < j} \frac{c_i c_j}{2m} + \sum_i \frac{c_i^2}{4m} = \sum_{ij} \frac{c_i c_j}{4m} = m. \quad (12.10)$$

We can also calculate the average number of ends of edges connected to a node i , i.e., its average degree $\langle k_i \rangle$. Allowing for the fact that a self-edge contributes two ends of edges to the degree, we get

$$\langle k_i \rangle = 2p_{ii} + \sum_{j(\neq i)} p_{ij} = \frac{c_i^2}{2m} + \sum_{j(\neq i)} \frac{c_i c_j}{2m} = \sum_j \frac{c_i c_j}{2m} = c_i. \quad (12.11)$$

In other words, the parameters c_i appearing in the definition of p_{ij} , Eq. (12.8), are the expected degrees in this model, just as the parameter c in $G(n, p)$ is the average degree of a node. The *actual* degree of a node could in principle take any value, depending on the luck of the draw about which edges happen to get randomly created and which do not. In general, the degree of node i will have a Poisson distribution with mean c_i , meaning that it will usually be quite

⁴As before, p_{ij} should really be regarded as the expected number of edges between i and j rather than the probability and in fact the proper formulation of the model is that we place a Poisson-distributed number of edges with mean p_{ij} between each pair of nodes i, j . Thus the model can in principle have multiedges as well as self-edges, just as in the configuration model. In the limit of large m and constant c_i , however, the probability and the expected number again become equal, and the density of multiedges tends to zero, so the distinction is unimportant.

⁵Another way of putting this is that the average value $\langle A_{ij} \rangle$ of an element of the adjacency matrix is simply $\langle A_{ij} \rangle = c_i c_j / 2m$ for all i, j —recall that the diagonal element A_{ii} of the adjacency matrix is defined to be twice the number of self-edges at node i , and this compensates for the extra factor of two in Eq. (12.8).

narrowly distributed about c_i , but there will always be some spread around this value, unless c_i is zero.⁶ Note that c_i does not have to be an integer, unlike the degrees k_i appearing in the configuration model.

This model is sometimes called the Chung–Lu model, after two of the first researchers to study it in detail [103]. In many ways it is an easier model to work with than the configuration model for the same reason that $G(n, p)$ is easier to work with than $G(n, m)$: the edges are independent random variables. On the other hand it has the disadvantage that we only get to specify the expected number of edges m and the expected degrees c_i in the network and not the actual values. This in turn means that we cannot choose the exact degree *distribution* of our network,⁷ and because the degree distribution plays such an important role in the study of networks this has made the model less attractive to researchers, despite its other advantages. Most calculations, as a result, are made using the configuration model and this is the direction we will take in this book as well. The approach behind the Chung–Lu model, however—if not the exact model itself—does have an important application in the degree-corrected stochastic block model, a model of networks with community structure that we study in Sections 12.11.6 and 14.4.1.

12.2 EXCESS DEGREE DISTRIBUTION

Having defined the configuration model, we now turn to a study of its properties. We begin our discussion with some fundamental observations about the model—and networks in general—that will prove central to later developments.

Consider a configuration model⁸ with degree distribution p_k , meaning that a fraction p_k of the nodes have degree k . Equivalently, p_k is the probability that a node chosen uniformly at random from the network has degree k . But suppose instead that we take a node (randomly chosen or not) and follow one of its edges (assuming it has at least one) to the node at the other end. What is the probability that this node will have degree k ?

⁶That the degree has a Poisson distribution is a standard result from probability theory: the degree is, in this case, the sum of a set of Poisson random variables representing the edges, and any sum of Poisson variables itself follows a Poisson distribution.

⁷It is easy to see that there are some degree distributions that the model cannot reproduce at all—any distribution for which p_k is exactly zero for any k , for instance, since there is always a non-zero probability that a node can have degree k for any value of k .

⁸We can consider either the standard version of the model in which the degree sequence is fixed, as in Section 12.1, or the version in which the degrees are drawn at random from the distribution p_k . The results will not depend on which we use.

The answer cannot just be p_k . For instance, there is no way to reach a node with degree zero by following an edge in this way, because a node with degree zero has no edges. So the probability of reaching a node of degree zero is itself zero, and not p_0 .

In fact, the correct probability for general k is not hard to calculate. We know that an edge emerging from a node in a configuration model network has equal chance of terminating at any edge “stub” anywhere else in the network (see Section 12.1). Since there are $\sum_i k_i = 2m$ stubs in total, or $2m - 1$ excluding the one at the beginning of our edge, and k of them are attached to any particular node with degree k , our edge has probability $k/(2m - 1)$ of ending at any particular node of degree k . In the limit where m becomes large we can ignore the -1 and just write this as $k/2m$.

Given that p_k is the total fraction of nodes in the network with degree k , the total number of such nodes is np_k , and hence the probability of our edge attaching to *any* node with degree k is

$$\frac{k}{2m} \times np_k = \frac{kp_k}{\langle k \rangle}, \quad (12.12)$$

where $\langle k \rangle$ is the average degree over the whole network and we have made use of the fact that $\langle k \rangle = 2m/n$, Eq. (6.15).

Thus the probability that we reach a node of degree k upon following an edge in this way is proportional not to p_k but to kp_k . To put that another way, the node you reach by following an edge is not a typical node in the network. It is more likely to have high degree than a typical node. Physically, the reasoning behind this observation is that a node with degree k has k edges attached to it, and you can reach that node by following any one of them. Thus if you choose an edge and follow it you have k times the chance of reaching a node with degree k than you have of reaching a node with degree 1.

It is important to recognize that this calculation is specific to the configuration model. It relies on the fact that edges in the configuration model are equally likely to end at any stub in the network. In the real world, this is not true: the degrees of adjacent nodes in networks are often correlated (see Section 7.7), which means the probability of reaching a node of degree k when we follow an edge depends on what node we are coming from.⁹ Nonetheless, the basic result is found to apply approximately to many real-world networks,

⁹On the other hand, if we pick a random *edge* in a network and follow it to one of its ends, then the degree of the node we reach is distributed according to (12.12), regardless of whether degrees are correlated. Picking random edges in networks is, however, not a very realistic exercise—there are few real-world processes that are equivalent to picking random edges.

which is one of the reasons why insights gained from the configuration model are useful for understanding the world around us.

Equation (12.12) has some strange and counterintuitive consequences. As an example, consider a randomly chosen node in the configuration model and let us calculate the average degree of a neighbor of that node. If we were using the configuration model to model a friendship network, for instance, the average degree of a network neighbor would correspond to the average number of friends your friend has. This number is the average of k over the probability in Eq. (12.12), which we get by multiplying the probability by k and then summing:

$$\text{average degree of a neighbor} = \sum_k k \frac{k p_k}{\langle k \rangle} = \frac{\langle k^2 \rangle}{\langle k \rangle}. \quad (12.13)$$

Note that the average degree of a neighbor is thus different from the average degree $\langle k \rangle$ of a typical node in the network. In fact, it is in general larger, as we can show by calculating the difference

$$\frac{\langle k^2 \rangle}{\langle k \rangle} - \langle k \rangle = \frac{1}{\langle k \rangle} (\langle k^2 \rangle - \langle k \rangle^2) = \frac{\sigma_k^2}{\langle k \rangle}, \quad (12.14)$$

where $\sigma_k^2 = \langle k^2 \rangle - \langle k \rangle^2$ is the variance of the degree distribution. The variance, which is the square of the standard deviation, is necessarily non-negative and indeed is strictly positive unless every single node in the network has the same degree. Let us assume that there is some variation in the degrees so that σ_k^2 is greater than zero. The average degree $\langle k \rangle$ is also greater than zero, unless all nodes have degree zero. Thus $\sigma_k^2 / \langle k \rangle > 0$ and Eq. (12.14) implies that $\langle k^2 \rangle / \langle k \rangle - \langle k \rangle > 0$, or

$$\frac{\langle k^2 \rangle}{\langle k \rangle} > \langle k \rangle. \quad (12.15)$$

In other words, the average degree of the neighbor of a node is greater than the average degree of a node. In colloquial terms, “Your friends have more friends than you do.”

This result is known as the *friendship paradox*, and at first sight it appears very strange. Certainly it seems likely that there will be some nodes in the network with higher degree than the average. But there will also be some that have lower degree and when you average over all neighbors of all nodes surely the two should cancel out? Shouldn’t the average degree of a neighbor be the same as the average degree in the network as a whole? Yet Eq. (12.15) tells us that this is not so, and it really is correct. You can create a configuration model network on a computer and compare the average degrees of neighbors

The ratio $\langle k^2 \rangle / \langle k \rangle$ that appears here crops up repeatedly in the study of networks. It will appear, for instance, in Chapter 15 when we study percolation theory and in Chapter 16 when we study the spread of disease, as well as several more times in this chapter.

to the average degree in the network as a whole and confirm that the former is indeed always larger than the latter. Even more remarkably, as first shown by Feld [170], you can do the same thing with real networks and, although the configuration model results don't apply exactly to these networks, the basic principle still seems to hold. Here, for instance, are some measurements for two academic coauthorship networks and a snapshot of the structure of the Internet at the autonomous system level:

Network	n	Average degree	Average neighbor degree	$\frac{\langle k^2 \rangle}{\langle k \rangle}$
Biologists	1 520 252	15.5	68.4	130.2
Mathematicians	253 339	3.9	9.5	13.2
Internet	22 963	4.2	224.3	261.5

According to these results a biologist's collaborators have, on average, more than four times as many collaborators as they do themselves. On the Internet, a node's neighbors have more than 50 times the average degree! Note that in each case the configuration model value of $\langle k^2 \rangle / \langle k \rangle$ overestimates the real average neighbor degree, in some cases by a substantial margin. This is typical of calculations using simplified network models like the configuration model: they can give you a general feel for the kind of behavior one might expect to see, but they usually don't give quantitatively accurate predictions for real networks.

The fundamental reason for the friendship paradox is that when you go through the nodes of a network and average the degrees of the neighbors of each one, many of those neighbors will appear in more than one average. In fact, a node with degree k will appear as a neighbor of exactly k other nodes, and hence appear in k of the averages. This means that high-degree nodes are over-represented in the calculations compared with low-degree ones and it is this bias that pushes up the overall average value.

In most of the calculations that follow, we will be interested not in the total degree of the node at the end of an edge but in the number of edges attached to that node *other* than the one we arrived along. This number is called the *excess degree* of the node and it is just equal to the total degree minus one.¹⁰ We can calculate the probability distribution of the excess degree from Eq. (12.12). The probability q_k of having excess degree k is equal to the probability of having

¹⁰The excess degree cannot, however, be negative. Since the node at the end of an edge always has degree at least 1 (because of the edge we followed to reach it) the minimum value of the excess degree is zero.

There is no reason in principle why the configuration model should always overestimate. In some cases it could underestimate too.

total degree $k + 1$ and, putting $k \rightarrow k + 1$ in Eq. (12.12), we get

$$q_k = \frac{(k+1)p_{k+1}}{\langle k \rangle}. \quad (12.16)$$

(Note that the denominator is still just $\langle k \rangle$, and not $\langle k + 1 \rangle$, as you can verify for yourself by checking that Eq. (12.16) is correctly normalized so that $\sum_{k=0}^{\infty} q_k = 1$.)

The distribution q_k is called the *excess degree distribution* and it will play an important role in many of the calculations that follow. It is the probability distribution, for a node reached by following an edge, of the number of other edges attached to that node.

12.3 CLUSTERING COEFFICIENT

As a simple application of the excess degree distribution, let us calculate the clustering coefficient for the configuration model. Recall that the clustering coefficient is the average probability that two neighbors of a node are also neighbors of each other.

Consider a node v that has at least two neighbors, which we will denote i and j . Being neighbors of v , i and j are both at the other ends of edges from v , and hence the number of other edges connected to them, which we will denote k_i and k_j , are distributed according to the excess degree distribution of Eq. (12.16). The probability of an edge between i and j is then $k_i k_j / 2m$ (see Eq. (12.2)) and, averaging both k_i and k_j over the distribution q_k , we get an expression for the clustering coefficient thus:

$$\begin{aligned} C &= \sum_{k_i, k_j=0}^{\infty} q_{k_i} q_{k_j} \frac{k_i k_j}{2m} = \frac{1}{2m} \left[\sum_{k=0}^{\infty} k q_k \right]^2 \\ &= \frac{1}{2m \langle k \rangle^2} \left[\sum_{k=0}^{\infty} k(k+1) p_{k+1} \right]^2 = \frac{1}{2m \langle k \rangle^2} \left[\sum_{k=0}^{\infty} (k-1) k p_k \right]^2 \\ &= \frac{1}{n} \frac{[\langle k^2 \rangle - \langle k \rangle]^2}{\langle k \rangle^3}, \end{aligned} \quad (12.17)$$

where we have made use of $2m = n \langle k \rangle$, Eq. (6.15).

Like the clustering coefficient of the Poisson random graph, Eq. (11.11), this expression goes as n^{-1} and so vanishes in the limit of large system size, as long as the moments $\langle k \rangle$ and $\langle k^2 \rangle$ of the degree distribution are fixed. Hence, like the Poisson random graph, the configuration model appears to be an unpromising model for real-world networks with high clustering. Note, however,

that Eq. (12.17) contains the second moment $\langle k^2 \rangle$ in its numerator, which can become large, for instance in networks with power-law degree distributions (see Section 10.4.2). This can result in surprisingly large values of C in the configuration model, as discussed in Section 10.6.

12.4 LOCALLY TREE-LIKE NETWORKS

Another important property of the configuration model is that the networks it generates are *locally tree-like*, meaning that any local neighborhood in such a network takes the form of a tree. More precisely, if you start at any node in the network and form the set of all nodes at distance d or less from that starting node, the set will, with probability tending to 1 as $n \rightarrow \infty$, take the form of a tree. The proof of this fact follows the same lines as the proof in Section 11.6 that the small components in a random graph are trees.

Consider a large network, choose a starting node, and form the neighborhood consisting of nodes at distance d away or less. Suppose this neighborhood has s nodes in total. Since it is, by definition, connected together, it must contain at least $s - 1$ edges—the minimum number needed to connect it. If it has exactly $s - 1$ edges then it is a tree. If it has more then it is not a tree.

But the probability that it has more than $s - 1$ edges is very small. Each additional edge we add between two nodes i, j in the neighborhood has probability $k_i k_j / 2m$, as shown in Section 12.1.1. If the average values of degrees k_i, k_j within the neighborhood remain fixed as the size of the network—and hence m —becomes large, then this probability vanishes and the neighborhood will take the form of a tree.

Note, however, that the nodes in the neighborhood are all reached by following edges. We work outward from our starting node, following edges repeatedly to reach them, so their degrees are distributed not according to the normal degree distribution but according to the excess degree distribution, and hence their average is the average of the excess degree distribution. From Eq. (12.16), the average excess degree is

$$\sum_{k=0}^{\infty} k q_k = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} k(k+1)p_{k+1} = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} (k-1)k p_k = \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle}. \quad (12.18)$$

Thus the network will be locally tree-like provided this average remains constant as the network becomes large.¹¹

¹¹Technically, we only require that it grows slower than \sqrt{m} for the proof to work, but in all the cases we will consider it is constant.

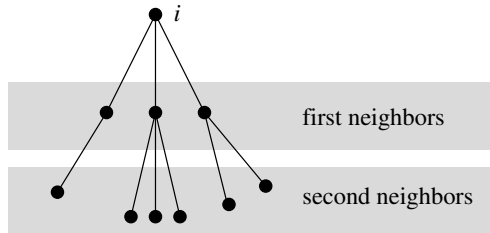


Figure 12.2: Calculation of the number of second neighbors of a node. The number of second neighbors of a node (top) is equal to the sum of the excess degrees of the first neighbors.

We can see the result of Section 12.3, that the clustering coefficient tends to zero as the network becomes large, as a special case of this more general result. If the network is locally tree-like, then it follows that two neighbors of a node cannot be connected, since if they were there would be a loop in the network and it would not be tree-like. Hence the clustering coefficient must go to zero.

The property that configuration model networks are locally tree-like will be crucial for many of the calculations we do in the remainder of this chapter.

12.5 NUMBER OF SECOND NEIGHBORS OF A NODE

An example of the application of the locally tree-like property arises when we look at the number of second neighbors of a node. The number of second neighbors varies from node to node in general, but we can calculate its average straightforwardly. Because local neighborhoods in the network take the form of trees, the number of second neighbors of a node i is simply equal to the sum of the excess degrees of the first neighbors—see Fig. 12.2. Thus, the average number of second neighbors is just the average excess degree multiplied by the number of first neighbors. The average excess degree is given by Eq. (12.18), so the average number of second neighbors of node i is $k_i(\langle k^2 \rangle - \langle k \rangle)/\langle k \rangle$, where k_i is i 's degree. If we now average over all nodes i , then k_i is replaced by its average $\langle k \rangle$, and number of second neighbors averaged over the whole network, which we denote c_2 , is

$$c_2 = \langle k \rangle \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} = \langle k^2 \rangle - \langle k \rangle. \quad (12.19)$$

In terms of this quantity, we can conveniently write the average excess degree of Eq. (12.18) as

$$\text{average excess degree} = \frac{c_2}{c_1}, \quad (12.20)$$

where we have introduced the notation $c_1 = \langle k \rangle$. (We previously called this quantity c , but we use c_1 here to emphasize the distinction between the numbers of first and second neighbors.)

We can take this approach further and calculate the mean number c_3 of neighbors at distance 3. The number of third neighbors is the sum of the excess degrees of the second neighbors, so its average is the average number of second neighbors times the average excess degree, or

$$c_3 = c_2 \times \frac{c_2}{c_1} = \frac{c_2^2}{c_1}. \quad (12.21)$$

Generalizing the argument, we see that every time we go one step further away from a node, the average number of neighbors at that distance increases by a factor of the average excess degree c_2/c_1 . Since the average number of neighbors at distance 1 is c_1 by definition, it follows that the general expression for the average number at distance d is

$$c_d = \left(\frac{c_2}{c_1}\right)^{d-1} c_1. \quad (12.22)$$

It's easy to confirm that this gives the correct expressions for the cases of $d = 1, 2$, and 3 .

12.6 GIANT COMPONENT

In the previous section, we showed that the average number of neighbors at distance d from a node in a configuration model network follows Eq. (12.22), which in turn implies that it either grows or falls off exponentially, depending on whether the ratio c_2/c_1 of the average numbers of first and second neighbors is greater or less than 1. This observation is strongly reminiscent of the argument we made in Section 11.5 for the appearance of a giant component in a random graph. There we argued that if the number of nodes you can reach within a certain distance is increasing with that distance (on average), then you must have a giant component in the network, while if it is decreasing there can be no giant component. Applying the same reasoning here, we conclude that the configuration model has a giant component if and only if we have

$$\frac{c_2}{c_1} > 1. \quad (12.23)$$

If this condition is not met and there is no giant component, the network can consist only of small components. Using Eq. (12.19) for c_2 and putting $c_1 = \langle k \rangle$, we can also write the condition as $\langle k^2 \rangle - \langle k \rangle > \langle k \rangle$ or

$$\langle k^2 \rangle - 2\langle k \rangle > 0. \quad (12.24)$$

This condition for the existence of a giant component in the configuration model was first given by Molloy and Reed [337] in 1995.¹²

Writing $\langle k \rangle = n^{-1} \sum_i k_i$ and $\langle k^2 \rangle = n^{-1} \sum_i k_i^2$, we can also express Eq. (12.24) as

$$\sum_i k_i(k_i - 2) > 0. \quad (12.25)$$

We note an interesting fact about this equation, that nodes of degree zero and degree two make no contribution to the sum, since terms in which $k_i = 0$ or $k_i = 2$ vanish. Thus we can add as many nodes of degree zero or two to the network as we like (or take them away) and it will make no difference to the existence or not of a giant component. We will see an example of this phenomenon in Section 12.6.1.

We can also calculate the size of the giant component, if there is one. The calculation is reminiscent of that for the Poisson random graph in Section 11.5. Consider any node in the network, choose one of its edges (assuming it has at least one), and follow that edge to the node at its other end. Let us define u to be the probability that this node does *not* belong to the giant component. Because, as we have seen, every edge in the configuration model is equally likely to attach to any “stub” in the network, this probability is the same no matter what node we start at or what edge we follow—the likelihood of landing in the giant component is the same in all cases.

To belong to the giant component, a node must be connected to the giant component via at least one of its neighbors. Or equivalently, a node does not belong to the giant component if (and only if) it is not connected to the giant component via any of its neighbors, which happens with probability u^k if it has k neighbors.

¹²This result has an interesting history. In the 1940s Flory [182] considered a model of branching polymers in which elemental units with a fixed number of “legs”—nodes with uniform degree, in effect—joined together to form connected clumps. He showed that, if the system was restricted to forming only trees, then there was a transition at which the polymer “gelled” to create a clump of units which corresponds to our giant cluster. In effect, Flory’s results were a special case of the solution given here for the uniform degree distribution, although they were not expressed in the language of networks. It was not until much later that Molloy and Reed, who appear to have been unaware of Flory’s work, gave the full solution for general degree distribution.

This assumes that the probabilities are all independent, which they will not be if, for instance, there are any direct connections between the neighbors—if two neighbors are connected by an edge, then one of them is in the giant component whenever the other is and their probabilities are not independent. In the limit of large network size, however, there are no such connections. The property of being locally tree-like ensures that there are no direct connections between neighbors, and not even any indirect ones along paths that run through other nodes (no matter how long those paths might be). So our assumption of independence is a good one and the probability of not being connected to the giant component is correctly given by u^k .

The average probability, over the entire network, that a node does not belong to the giant component is now given by the average of u^k over all values of the degree k , which is $\sum_k p_k u^k$, where p_k is, as usual, the degree distribution of the network. This particular sum arises often in the study of the configuration model (and other network models), so we give it its own notation:

$$g_0(u) = \sum_k p_k u^k. \quad (12.26)$$

The function $g_0(u)$ is called the *probability generating function* for the probability distribution p_k .

But if $g_0(u)$ is the average probability that a node does not belong to the giant component then

$$S = 1 - g_0(u) \quad (12.27)$$

is the probability that it *does* belong to the giant component, or equivalently the fraction of the network occupied by the giant component.

To make use of this result we still need to know the value of u , the average probability that a node is not connected to the giant component via one of its neighbors. This we calculate as follows. The probability that a node is not connected to the giant component via a particular neighboring node is equal to the probability that *that* node is not connected to the giant component via any of its other neighbors. If there are k of those other neighbors, then that probability is again u^k . But because we are talking about a neighboring node, k is now distributed according to the excess degree distribution q_k of Eq. (12.16), so the average probability of not being connected to the giant component via one's neighbor is $\sum_k q_k u^k$. But this probability is, by definition, just u , so we get a self-consistent equation for u thus:

$$u = \sum_k q_k u^k. \quad (12.28)$$

The sum in this equation also occurs often in the study of networks, so we define

$$g_1(u) = \sum_k q_k u^k. \quad (12.29)$$

This function is the probability generating function for the excess degree distribution. In terms of this generating function Eq. (12.28) reads

$$u = g_1(u). \quad (12.30)$$

Taken together, Eqs. (12.27) and (12.30), along with the definitions of the two generating functions, give us our solution for the size S of the giant component. Given the degree distribution and excess degree distribution of the network, we can calculate the generating functions from their definitions, Eqs. (12.26) and (12.29), then the value of u is given by the solution of (12.30) and, substituting this value into (12.27), we get our solution for S .

Although it is convenient to have separate notations for the two generating functions $g_0(u)$ and $g_1(u)$, they are not really independent since the excess degree distribution is itself defined in terms of the ordinary degree distribution via Eq. (12.16). Using Eq. (12.16) we can write $g_1(u)$ as

$$g_1(u) = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} (k+1)p_{k+1}u^k = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} kp_k u^{k-1} = \frac{1}{\langle k \rangle} g'_0(u), \quad (12.31)$$

where we have made use of (12.26) and g'_0 denotes the first derivative of g_0 with respect to its argument. But note also that

$$g'_0(1) = \sum_{k=0}^{\infty} kp_k = \langle k \rangle, \quad (12.32)$$

so

$$g_1(u) = \frac{g'_0(u)}{g'_0(1)}. \quad (12.33)$$

This convenient formula allows us to calculate $g_1(u)$ directly from $g_0(u)$, without needing to calculate the excess degree distribution itself.

12.6.1 EXAMPLE

Let us take a look at a concrete example and see how these formulas work out in practice. Consider, for instance, a network that has nodes only of degree 0, 1, 2, and 3, and no nodes of any higher degree. For such a network one only

has to specify the values of p_0, p_1, p_2 , and p_3 , the rest of the p_k being zero, and the generating functions $g_0(u)$ and $g_1(u)$ take the form

$$g_0(u) = p_0 + p_1u + p_2u^2 + p_3u^3, \quad (12.34)$$

$$g_1(u) = \frac{g'_0(u)}{g'_0(1)} = \frac{p_1 + 2p_2u + 3p_3u^2}{p_1 + 2p_2 + 3p_3} = q_0 + q_1u + q_2u^2. \quad (12.35)$$

Equation (12.30) is thus a quadratic equation in this case, $u = q_0 + q_1u + q_2u^2$, which has the solutions

$$u = \frac{1 - q_1 \pm \sqrt{(1 - q_1)^2 - 4q_0q_2}}{2q_2}. \quad (12.36)$$

However, since q_k is a probability distribution it must sum to unity, $q_0 + q_1 + q_2 = 1$, and hence $1 - q_1 = q_0 + q_2$. Using this result to eliminate q_1 from (12.36) we get

$$\begin{aligned} u &= \frac{(q_0 + q_2) \pm \sqrt{(q_0 + q_2)^2 - 4q_0q_2}}{2q_2} \\ &= \frac{(q_0 + q_2) \pm (q_0 - q_2)}{2q_2} \\ &= 1 \quad \text{or} \quad \frac{q_0}{q_2}. \end{aligned} \quad (12.37)$$

Recall that u represents the average probability that a node is not connected to the giant component via its neighbor. The solution $u = 1$ thus corresponds to a situation where there is no giant component in the network—no node belongs to the giant component because the probability of being connected to it is zero. But we also have a second solution $u = q_0/q_2$ which *can* give us a giant component. Which of these two solutions we use depends on whether or not the network does in fact contain a giant component.

In Section 12.6 we showed that there is a giant component if and only if the ratio of the average numbers of first and second neighbors of a node satisfies $c_2/c_1 > 1$. But recall that c_2/c_1 is also the average excess degree (Eq. (12.20)), which we can calculate directly for our example network:

$$\text{average excess degree} = \sum_{k=0}^{\infty} kq_k = q_1 + 2q_2 = 1 - q_0 + q_2. \quad (12.38)$$

This is greater than 1 if $q_2 - q_0 > 0$ or equivalently if $q_0/q_2 < 1$.

Thus we arrive at a simple and elegant conclusion regarding the two solutions, Eq. (12.37), for u . If the second solution q_0/q_2 is greater than or equal

to 1, then there is no giant component, meaning the first solution $u = 1$ applies (which is a good thing, since u is a probability and cannot be greater than 1 anyway, so the second solution would not be allowed). However, if the second solution is less than 1 then there is a giant component and we should adopt the second solution not the first for u :

$$u = \frac{q_0}{q_2} = \frac{p_1}{3p_3}, \quad (12.39)$$

where we have extracted the values of q_0 and q_2 from Eq. (12.35). This solution for u is less than 1 when

$$p_3 > \frac{1}{3}p_1. \quad (12.40)$$

In other words, there is a giant component if the number of nodes of degree three exceeds one-third the number of degree one. Note that the number of nodes of degree zero and degree two don't enter into the picture at all (except to the extent that their absence makes room for more nodes of the other degrees). As discussed in Section 12.6, this is a general result: nodes of degree zero and two never make any difference to the presence or absence of a giant component.

The size of the giant component in this example is given by Eq. (12.27) to be

$$S = 1 - g_0(u) = 1 - p_0 - \frac{p_1^2}{3p_3} - \frac{p_1^2 p_2}{9p_3^2} - \frac{p_1^3}{27p_3^3}. \quad (12.41)$$

Thus the size of the giant component does depend on p_0 and p_2 , even though its presence or absence does not.

12.6.2 GENERAL SOLUTION FOR THE SIZE OF THE GIANT COMPONENT

The example of the last section is unusual in that we can solve Eq. (12.30) exactly for the crucial parameter u . In most other cases exact solutions are not possible, but we can nonetheless get a good idea of the behavior of u as follows.

First, we note that the definition of the generating function $g_1(u)$ implies that $g_1(1) = \sum_k q_k = 1$, since q_k is a properly normalized probability distribution. Hence the equation $u = g_1(u)$, Eq. (12.30), always has a trivial solution $u = 1$, no matter what the degree distribution is. We saw a special case of this in the last section, but it is true in general for all configuration models. As previously, $u = 1$ corresponds to the situation where there is no giant component. But as we have seen, Eq. (12.30) can also have another solution that does give a giant component.

The function $g_1(u)$ is defined as a power series, Eq. (12.29), with coefficients equal to the probabilities q_k and hence all non-negative. That means that

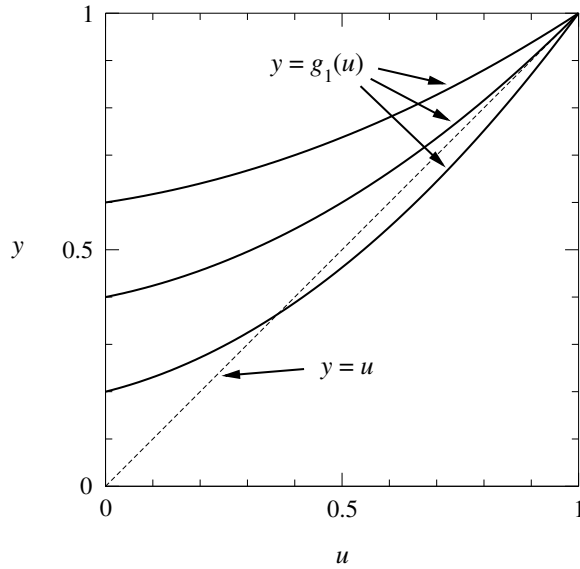


Figure 12.3: Graphical solution of Eq. (12.30). The solution of the equation $u = g_1(u)$ is given by the point at which the curve $y = g_1(u)$ intercepts the line $y = u$.

the derivatives of $g_1(u)$ are also all non-negative whenever $u \geq 0$ and hence that $g_1(u)$ is in general positive, an increasing function of its argument, and upward concave. Given that it takes the value 1 when $u = 1$, it must thus look qualitatively like one of the curves in Fig. 12.3. The solution of the equation $u = g_1(u)$ is then given by the intercept of the curve $y = g_1(u)$ with the line $y = u$ (the dotted line in the figure).

The trivial solution at $u = 1$ appears at the upper right in the figure and is always present, but there may or may not be another solution with $u < 1$ if the curve takes the right form. In particular, as the figure shows, we have a non-trivial solution at $u < 1$ if the slope $g'_1(1)$ of the curve at $u = 1$ is greater than the slope of the dotted line, that is, if

$$g'_1(1) > 1. \tag{12.42}$$

But, using Eq. (12.29) for $g_1(u)$ again, we have

$$g'_1(1) = \sum_{k=0}^{\infty} k q_k, \tag{12.43}$$

which is none other than the average excess degree again. So there is a solution

with $u < 1$ if and only if the average excess degree is greater than 1, which, as we showed in Section 12.6, is precisely the regime in which the network has a giant component.

In other words, the story is basically the same as it was for the example of the previous section. When there is no giant component, $u = 1$ is the only solution to Eq. (12.30) in the interval from zero to one, but when there is a giant component a second solution $u < 1$ appears and it is this solution that gives us the size of our giant component.

Although we cannot always solve for u exactly, we can calculate it numerically, on a computer. A simple strategy is to choose a starting value for u (the value $u = \frac{1}{2}$ works fine) and iterate Eq. (12.30), repeatedly feeding u in on the right and getting a new value out on the left until the result converges. Fifty iterations are usually enough to give an accurate figure. Once we have the value of u , we can substitute it into Eq. (12.27) to get the fraction of the network occupied by the giant component.

12.7 SMALL COMPONENTS

Having looked in some detail at the behavior of the giant component in the configuration model, let us turn to the small components. As in with the Poisson random graph of Chapter 11, the small components can in principle take a range of different sizes, but it is relatively straightforward to calculate their average size. The calculation follows lines similar to those of the equivalent calculation for the Poisson random graph (Section 11.6), though the arguments are somewhat more involved, as we will see.

As before we focus on the average size of the component to which a randomly chosen small-component node i belongs.¹³ The situation is depicted in Fig. 12.4 (which is actually the same figure we previously used for the Poisson random graph on page 357, but it works just as well as an illustration of the configuration model). As we argued in Section 12.4, every local neighborhood of a node in the configuration model, no matter how large, is a tree in the large- n limit, and hence all small components (which are the neighborhoods of their constituent nodes) must be trees, just as in the Poisson random graph. But in

¹³Note that this is not the same as just the average size of a small component. As noted previously in footnote 7 on page 359, there are more nodes in a larger component than in a smaller one, which means that the average size of the component to which a node belongs is biased toward larger components. Nonetheless, this average turns out to be the most useful measure of component size in most practical situations, and (following most other work in this area) it is the one we use in this book.

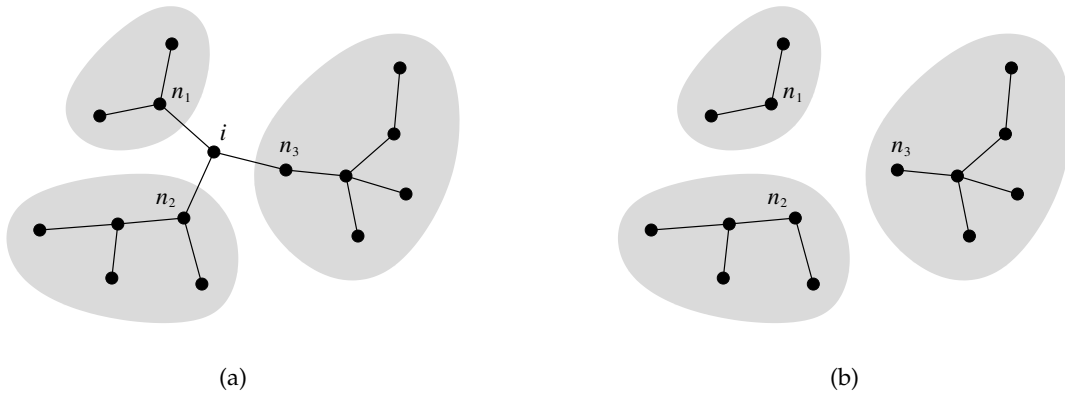


Figure 12.4: The size of one of the small components in the configuration model. (a) The size of the component to which a node i belongs is the sum of the number of nodes in each of the subcomponents (shaded regions) reachable via i 's neighbors n_1, n_2, n_3 , plus one for i itself. (b) If node i is removed the subcomponents become components in their own right.

that case the arguments of Section 11.6 again apply: the sets of nodes reachable along each of i 's edges (shaded regions in the figure) cannot be connected to one another, except via i , and hence the size of the component is simply the sum of the sizes of these sets, plus 1 for node i itself.

If node i has degree k and the sizes of the sets are $t_1 \dots t_k$ then, following the same argument as for Eq. (11.26), the average size of the component it belongs to is $1 + k\langle t \rangle$ and, averaging over all nodes i in small components,

$$\langle s \rangle = 1 + \langle k \rangle_{\text{small}} \langle t \rangle, \quad (12.44)$$

which is the same as Eq. (11.27) for the Poisson case. As there, $\langle k \rangle_{\text{small}}$ is the average degree of a node in a small component and $\langle t \rangle$ is the average size of the set of nodes reached by following an edge. To evaluate Eq. (12.44) we now need to find the values of these two quantities.

12.7.1 DEGREES OF NODES IN THE SMALL COMPONENTS

The distribution of degrees of nodes in the small components is different in general from that in the network as whole because nodes with higher degrees are less likely to be in the small components. The probability that a node belongs to a small component, given that it has degree k , is equal to the probability that it doesn't belong to the giant component, which, as described in Section 12.6, is

just u^k , where u is the solution of Eq. (12.30). We can use this result to calculate the probability that a node has degree k given that it is in a small component, by applying Bayes rule:

$$\begin{aligned} P(\text{degree } k | \text{small component}) \\ = P(\text{small component} | \text{degree } k) \frac{P(\text{degree } k)}{P(\text{small component})}. \end{aligned} \quad (12.45)$$

But $P(\text{degree } k) = p_k$ by definition and $P(\text{small component}) = 1 - S = g_0(u)$, where S is the size of the giant component, Eq. (12.27). So

$$P(\text{degree } k | \text{small component}) = \frac{p_k u^k}{g_0(u)}. \quad (12.46)$$

And the average degree of a node in a small component is just the average of this distribution:

$$\langle k \rangle_{\text{small}} = \frac{1}{g_0(u)} \sum_{k=0}^{\infty} k p_k u^k = \frac{u g_0'(u)}{g_0(u)}. \quad (12.47)$$

12.7.2 AVERAGE NUMBER OF NODES REACHED ALONG AN EDGE

To calculate the average size $\langle t \rangle$ of the sets of nodes in Fig. 12.4 we use the same trick that we used previously for the Poisson random graph in Section 11.6: we remove node i from the picture. When we do this, those sets of nodes become small components in their own right and then the sizes of the sets are given by the sizes of the components to which the neighbors n_1, n_2, \dots belong. Suppose one of these neighbors has degree k . As we showed in Section 12.7, the average size of the small component to which a node of degree k belongs is $1 + k\langle t \rangle$, and the same expression applies here too. The average size of a set is then given by the average of this expression over the degrees k of neighbors of nodes in the small components $\langle t \rangle = 1 + \langle k \rangle_{\text{neighbor}} \langle t \rangle$, or

$$\langle t \rangle = \frac{1}{1 - \langle k \rangle_{\text{neighbor}}}. \quad (12.48)$$

And what is the average degree $\langle k \rangle_{\text{neighbor}}$ of a neighbor? The neighbors are nodes we reach by following an edge from node i , so the “degree” k in this case is actually the excess degree. But here we must be careful: in the same way that the nodes in the small components don’t have the same degree distribution as nodes in the network as a whole, they don’t have the same excess degree distribution either.

As we showed in Section 12.2, the excess degree distribution in the network as a whole, Eq. (12.16), is proportional to k times the normal degree distribution p_k . By the same argument, the excess degree distribution in the small components is proportional to k times the degree distribution from Eq. (12.46), i.e., proportional to $kp_k u^k$. Normalizing to get a proper probability distribution, and noting that (as usual) the excess degree k is one less than the total degree, the excess degree distribution in the small components is

$$\frac{(k+1)p_{k+1}u^{k+1}}{\sum_k kp_k u^k} = \frac{(k+1)p_{k+1}u^{k+1}}{u g'_0(u)}. \quad (12.49)$$

Thus the average excess degree in the small components is

$$\begin{aligned} \langle k \rangle_{\text{neighbor}} &= \frac{1}{u g'_0(u)} \sum_{k=0}^{\infty} k(k+1)p_{k+1}u^{k+1} = \frac{1}{u g'_0(u)} \sum_{k=0}^{\infty} (k-1)kp_k u^k \\ &= \frac{u g''_0(u)}{g'_0(u)}. \end{aligned} \quad (12.50)$$

Putting everything together and combining Eqs. (12.44), (12.47), (12.48), and (12.50), we now have

$$\langle s \rangle = 1 + \frac{u g'_0(u)}{g_0(u)} \frac{1}{1 - u g''_0(u)/g'_0(u)}. \quad (12.51)$$

We can simplify this somewhat by making use of Eqs. (12.30) and (12.33) to write

$$g_1(u) = \frac{g'_0(u)}{g'_0(1)}, \quad g'_1(u) = g_1(u)g'_0(1) = u g'_0(1), \quad g''_1(u) = g'_1(u)g'_0(1), \quad (12.52)$$

so that

$$\langle s \rangle = 1 + \frac{u^2 g'_0(1)}{g_0(u)[1 - g'_1(u)]}. \quad (12.53)$$

This is the average size of the component to which a small-component node in the configuration model belongs.

Equation (12.53) is a relatively complicated expression, but a simple case occurs when we are in the region where there is no giant component. In this region we have $u = 1$ (see the discussion in Section 12.6.2) and hence

$$\langle s \rangle = 1 + \frac{g'_0(1)}{1 - g'_1(1)}, \quad (12.54)$$

where we have made use of the fact that $g_0(1) = \sum_k p_k = 1$. Thus, the average size of the component to which a node belongs diverges at the point where $g'_1(1) = 1$ —the point at which the curve in Fig. 12.3 is exactly tangent to the dotted line (the middle curve in the figure), which is, as discussed in Section 12.6.2, precisely the point at which the giant component first appears.

Thus the picture we have is similar to that shown in Fig. 11.5 for the Poisson random graph, in which the typical size of the component to which a node belongs grows larger and larger until we reach the point, the phase transition, where the giant component appears, at which it diverges. Beyond this point a giant component appears and grows larger, but the small components shrink in size again.

Equation (12.54) can also be expressed in a couple of other forms that may be useful in some circumstances. From Eq. (12.43) we know that $g'_1(1)$ is equal to the average excess degree in the network, which is also equal to the ratio c_2/c_1 of the average numbers of first and second neighbors of a node, Eq. (12.20). Meanwhile $g'_0(1) = \sum_k k p_k = \langle k \rangle = c_1$, so

$$\langle s \rangle = 1 + \frac{c_1^2}{c_1 - c_2}, \tag{12.55}$$

so that the average size of the component a node belongs to is dictated entirely by the mean numbers of first and second neighbors. Alternatively, we could write $c_2 = \langle k^2 \rangle - \langle k \rangle$ as in Eq. (12.19) and $c_1 = \langle k \rangle$ to get

$$\langle s \rangle = 1 + \frac{\langle k \rangle^2}{2\langle k \rangle - \langle k^2 \rangle}, \tag{12.56}$$

so $\langle s \rangle$ is also specified entirely by the first and second moments of the degree distribution. Equation (12.56) can be evaluated easily given only a knowledge of the degrees in the network, avoiding the need to calculate any generating functions.

12.8 NETWORKS WITH POWER-LAW DEGREE DISTRIBUTIONS

As we saw in Section 10.4, some networks have degree distributions that approximately obey a power law. As an example application of the machinery developed in this chapter, let us look at the properties of configuration model networks with power-law degree distributions.

There are various forms used to represent power laws in practice but the simplest choice is the “pure” power law introduced in Section 10.4.2:

$$p_k = \begin{cases} 0 & \text{for } k = 0, \\ k^{-\alpha} / \zeta(\alpha) & \text{for } k \geq 1, \end{cases} \tag{12.57}$$

where

$$\zeta(\alpha) = \sum_{k=1}^{\infty} k^{-\alpha} \quad (12.58)$$

is the Riemann zeta function—see Eq. (10.13) and the surrounding discussion. There is no closed-form expression for the zeta function, but there exist good numerical methods for calculating its value accurately, and many programming languages and numerical software packages include functions to calculate it, so Eq. (12.57) is reasonably convenient to work with in practice.

Using the results of the previous sections we can now say, for instance, whether there will be a giant component in a configuration model network with this degree distribution. Equation (12.24) tells us that there will be a giant component if and only if

$$\langle k^2 \rangle - 2\langle k \rangle > 0. \quad (12.59)$$

In the present case

$$\langle k \rangle = \sum_{k=0}^{\infty} k p_k = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} k^{-\alpha+1} = \frac{\zeta(\alpha-1)}{\zeta(\alpha)}, \quad (12.60)$$

and

$$\langle k^2 \rangle = \sum_{k=0}^{\infty} k^2 p_k = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} k^{-\alpha+2} = \frac{\zeta(\alpha-2)}{\zeta(\alpha)}. \quad (12.61)$$

Thus there is a giant component if

$$\zeta(\alpha-2) > 2\zeta(\alpha-1). \quad (12.62)$$

Figure 12.5 shows this inequality in graphical form. The two curves in the figure show the values of $\zeta(\alpha-2)$ and $2\zeta(\alpha-1)$ as functions of α and, as we can see, the inequality (12.62) is satisfied only for sufficiently low values of α , below the dotted line in the figure. A numerical solution of the equation $\zeta(\alpha-2) = 2\zeta(\alpha-1)$ gives a value of $\alpha = 3.4788\dots$ for the position of this line, so the network will have a giant component only when $\alpha < 3.4788$, a result first given by Aiello *et al.* in 2000 [10].

In practice this result is of somewhat limited utility because it applies only for the pure power law. In general, other distributions with power-law tails but different behavior for low k will have different thresholds at which the giant

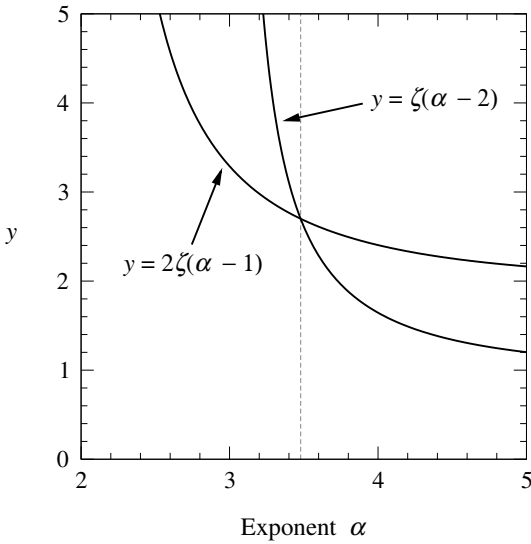


Figure 12.5: Graphical solution of Eq. (12.62). The configuration model with a pure power-law degree distribution, Eq. (12.57), has a giant component if $\zeta(\alpha-2) > 2\zeta(\alpha-1)$. This happens for values of α below the crossing point of the two curves, marked with the vertical dashed line.

component appears. There is, however, a general result we can derive that applies to all distributions with power-law tails. In Section 10.4.2 we noted that the second moment $\langle k^2 \rangle$ diverges for any distribution with a power-law tail with exponent $\alpha \leq 3$, while the first moment $\langle k \rangle$ remains finite so long as $\alpha > 2$. This means that Eq. (12.59) is always satisfied for any configuration model with a power-law tail to its degree distribution so long as α lies in the range $2 < \alpha \leq 3$, and hence there will always be a giant component no matter what else the distribution does. For $\alpha > 3$, on the other hand, there may or may not be a giant component, depending on the precise functional form of the degree distribution. (For $\alpha \leq 2$ it turns out that there is always a giant component, although more work is needed to demonstrate this.) Note that, as discussed in Section 10.4, most observed values of α for real-world networks lie in the range $2 < \alpha \leq 3$ and hence we tentatively expect such networks to have a giant component, although we must always bear in mind that the configuration model is a simplified model of a network and is not a foolproof guide to the behavior of real networks.

Returning to the pure power law let us calculate the size S of the giant component, when there is one. The generating function $g_0(u)$ for the degree distribution, Eq. (12.57), is

$$g_0(u) = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} k^{-\alpha} u^k. \quad (12.63)$$

The sum cannot be expressed in closed form, so we will just leave it as a sum for now. The generating function $g_1(u)$ for the excess degree distribution can be computed from Eq. (12.33), which gives:

$$g_1(u) = \frac{\sum_{k=1}^{\infty} k^{-(\alpha-1)} u^{k-1}}{\sum_{k=1}^{\infty} k^{-(\alpha-1)}} = \frac{1}{\zeta(\alpha-1)} \sum_{k=1}^{\infty} k^{-\alpha+1} u^{k-1}. \quad (12.64)$$

Now the crucial equation (12.30) for the probability u that a neighboring node is not in the giant component reads

$$u = \frac{1}{\zeta(\alpha-1)} \sum_{k=0}^{\infty} (k+1)^{-\alpha+1} u^k. \quad (12.65)$$

In general there is no closed-form solution for this equation, but we do note some interesting points. In particular, note that if $\zeta(\alpha-1)$ diverges we will get a solution $u = 0$. And indeed $\zeta(\alpha-1)$ does diverge, at $\alpha = 2$ and all values

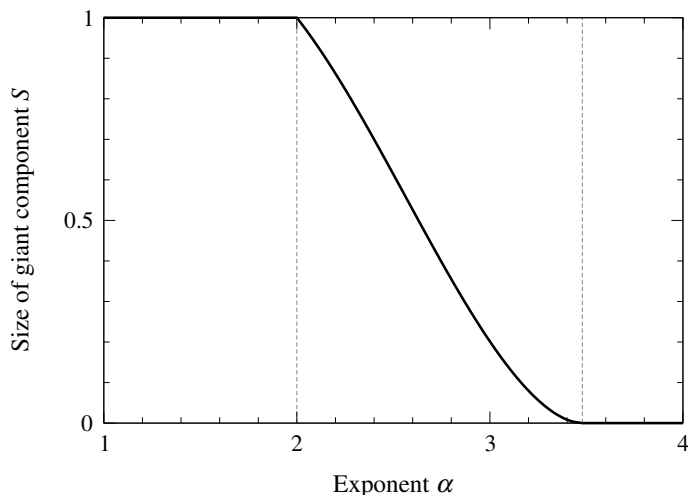


Figure 12.6: Size of the giant component for the configuration model with a power-law degree distribution. This plot shows the fraction of the network filled by the giant component as a function of the exponent α of the power law, calculated by numerical solution of Eqs. (12.27) and (12.65). The dotted lines mark the value $\alpha = 2$, below which the giant component has size 1, and the value $\alpha = 3.4788$, above which there is no giant component.

below, as one can readily verify from the definition in Eq. (12.58).¹⁴ Thus, for $\alpha \leq 2$ we have $u = 0$ and Eq. (12.27) then tells us that the giant component has size $S = 1 - g_0(0) = 1 - p_0$. However, for our particular choice of degree distribution, Eq. (12.57), there are no nodes with degree zero, and hence $p_0 = 0$ and $S = 1$. That is, the giant component fills the entire network and there are no small components at all!

Technically, this statement is not quite correct. There is always some chance that, for instance, a node of degree 1 will connect to another node of degree 1, forming a small component of size 2. What we have shown is that the probability that a randomly chosen node belongs to a small component is zero in the limit of large n , i.e., that the small components fill a fraction of the network that vanishes as $n \rightarrow \infty$. In the language used by mathematicians, a randomly chosen node belongs to the giant component “with high probability,” meaning

¹⁴Traditionally $\zeta(x)$ is actually defined to have finite values below $x = 1$ by analytic continuation, but in our case we are really interested in the value of the sum $\sum_{k=1}^{\infty} k^{-x}$, which diverges for all $x \leq 1$.

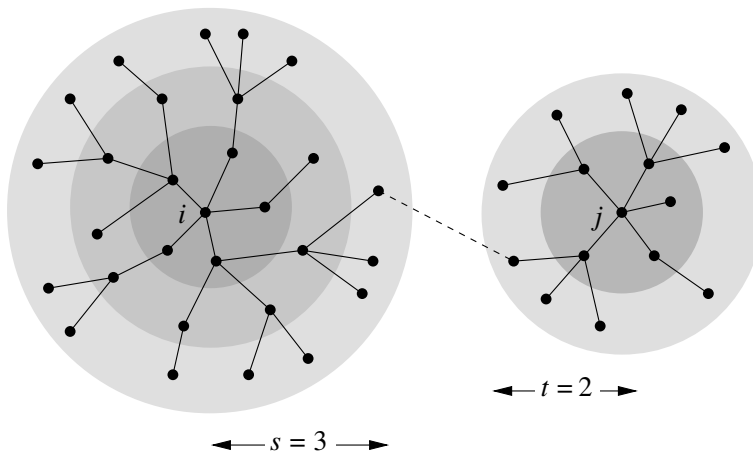


Figure 12.7: Neighborhoods of two nodes in a configuration model network. We consider the sets of nodes at distances s and t respectively from i and j . If there is an edge between any node in one set and any node in the other (dashed line), then there is a path between i and j of length $s + t + 1$.

it is technically possible to observe another outcome, but the probability is vanishingly small in the limit of large n .

Thus, our picture of the pure power-law configuration model is one in which there is a giant component for values of $\alpha < 3.4788$ and that giant component fills essentially the entire network when $\alpha \leq 2$. In the region between $\alpha = 2$ and $\alpha = 3.4788$ there is a giant component but it does not fill the whole network and some portion of the network consists of small components. If $\alpha > 3.4788$ there are only small components. To confirm this picture, Fig. 12.6 shows the size of the giant component extracted from a numerical solution of Eq. (12.65) using the method described at the end of Section 12.6. As we can see, it fits nicely with our expectations.

12.9 DIAMETER

We can also calculate the typical diameter of networks generated using the configuration model. The calculation is a variant of the one we used in Section 11.7 for the Poisson random graph. As in that case, we ask about the shortest path between two nodes i and j and consider two sets of nodes: those at distance s from i and those at distance t from j , as shown in Fig. 12.7 (which is the same figure we used in Chapter 11—it works equally well for the configuration

model). If there is a direct connection between the “surfaces” of these two sets of nodes, as indicated by the dashed line in the figure, then the shortest path from i to j has length no greater than $s + t + 1$. Conversely, if there is no such connection, then the shortest path must be longer than $s + t + 1$.

The probability of an edge between any individual pair of nodes u, v on the two surfaces is $k_u k_v / 2m$ following Eq. (12.2). But, since nodes on the surface are reached by following a sequence of edges from the starting points s and t , the relevant degrees k_u, k_v are the *excess* degrees of the corresponding nodes, which follow the excess degree distribution of Eq. (12.16). Averaging both k_u and k_v over this distribution, we get two factors of the average excess degree, which is c_2/c_1 (see Eq. (12.20)), where c_1 is the mean degree of a node in the network and c_2 is the mean number of second neighbors of a node, Eq. (12.19). Hence the average probability of an edge is $(c_2/c_1)^2/2m$.

The total number of pairs of nodes between which such an edge could fall, joining our two surfaces, is the product of the sizes of the surfaces. From Eq. (12.22) we know that these sizes are $(c_2/c_1)^{s-1} c_1$ and $(c_2/c_1)^{t-1} c_1$ on average, so the product of the sizes is

$$\left(\frac{c_2}{c_1}\right)^{s-1} c_1 \times \left(\frac{c_2}{c_1}\right)^{t-1} c_1 = \left(\frac{c_2}{c_1}\right)^{s+t-2} c_1^2. \quad (12.66)$$

And the probability that there is no connection between any of these pairs, which is also the probability that the distance d_{ij} between i and j is greater than $s + t + 1$, is:

$$P(d_{ij} > s + t + 1) = \left[1 - \frac{(c_2/c_1)^2}{2m}\right]^{(c_2/c_1)^{s+t-2} c_1^2}. \quad (12.67)$$

Taking the log of both sides and expanding in powers of $(c_2/c_1)^2/2m$ we get

$$\ln P(d_{ij} > s + t + 1) \simeq -\frac{c_1}{n} \left(\frac{c_2}{c_1}\right)^{s+t}, \quad (12.68)$$

where we have made use of $c_1 = 2m/n$ (Eq. (6.15)) and the approximate equality becomes exact in the limit of large network size, where both n and m diverge.

Taking exponentials again and defining $\ell = s + t + 1$, we get

$$P(d_{ij} > \ell) = \exp\left[-\frac{c_1}{n} \left(\frac{c_2}{c_1}\right)^{\ell-1}\right]. \quad (12.69)$$

By definition, the diameter of the network is that value of ℓ such that the probability that d_{ij} exceeds ℓ is zero. As n becomes large, Eq. (12.69) tells us that this happens only if $(c_2/c_1)^\ell$ grows faster than n . In other words

$(c_2/c_1)^\ell = an^{1+\epsilon}$ for some constant a and $\epsilon \rightarrow 0$ from above. Rearranging for ℓ , we then have

$$\ell = \frac{\ln a}{\ln(c_2/c_1)} + \lim_{\epsilon \rightarrow 0} \frac{(1 + \epsilon) \ln n}{\ln(c_2/c_1)} = A + \frac{\ln n}{\ln(c_2/c_1)}, \quad (12.70)$$

where A is another constant. The value of A is unknown, but in the limit of large n the second term in (12.70) dominates and we can ignore A , so to leading order the diameter of the configuration model is just $\ln n / \ln(c_2/c_1)$. This result, which was first derived by Chung and Lu [103], demonstrates among other things that networks drawn from the configuration model display the small-world effect in the sense of Section 11.7: path lengths between nodes in the network are of order $\ln n$ or shorter, and hence grow only very slowly with network size n . This gets around some of the objections we raised about the Poisson random graph in Section 11.8. That model too shows the small-world effect, but it is not a plausible model of most real networks for a range of reasons but particularly because of its unrealistic degree distribution. The configuration model can assume any degree distribution and still shows the small-world effect, making the result more believable.

12.10 GENERATING FUNCTION METHODS

In the previous sections of this chapter we have shown how to calculate a number of properties of configuration model networks using conventional methods of algebra and probability theory. We could continue in a similar vein, but the arguments—particularly in the past few sections—have already become quite involved and hard to follow. In practice, calculations for the configuration model, and many other network models as well, are often carried out using a set of more sophisticated mathematical techniques that manipulate generating functions. It takes some work to master these techniques, but the investment is worth it: once you have them under your belt, many network calculations become much more straightforward, and even advanced calculations that would have been daunting before become manageable. In this section, we introduce the general theory behind these generating function techniques, then demonstrate how it can be applied to a range of calculations for the configuration model.

12.10.1 GENERATING FUNCTIONS

The fundamental mathematical tool that we will use in the following sections is the probability generating function. We have already seen two examples

of generating functions earlier in the chapter, those for the degree and excess degree distributions, Eqs. (12.26) and (12.29). In this section we give a more formal introduction to generating functions and their properties. Readers interested in pursuing the mathematics of generating functions further may like to look at the book by Wilf [471].¹⁵

Suppose we have a probability distribution for a non-negative integer variable, such that p_k is the probability that the variable takes the value k . A good example of such a distribution is the distribution of the degrees of randomly chosen nodes in a network. If the fraction of nodes in a network with degree k is p_k then p_k is also the probability that a randomly chosen node in the network will have degree k .

The *generating function* for the probability distribution p_k is the polynomial

$$g(z) = p_0 + p_1z + p_2z^2 + p_3z^3 + \dots = \sum_{k=0}^{\infty} p_k z^k. \quad (12.71)$$

To be technically correct, this is a *probability generating function*, a name intended to distinguish it from another common type of function, the *exponential generating function*, which appears in certain types of counting problems. We will not use exponential generating functions in this book, so for us all generating functions will be probability generating functions.

If we know the generating function for a probability distribution p_k then we can recover the values of p_k by differentiating:

$$p_k = \frac{1}{k!} \left. \frac{d^k g}{dz^k} \right|_{z=0}. \quad (12.72)$$

Thus the generating function gives us complete information about the probability distribution and vice versa. We sometimes say that the probability distribution is *generated by* the function $g(z)$.

In effect, the distribution and the generating function are two different representations of the same thing. As we will see, it is easier in many cases to work with the generating function than with the probability distribution and this is the primary reason for their use in network calculations.

12.10.2 EXAMPLES

Right away let us look at some examples of generating functions. Taking a page from Section 12.6.1, suppose our variable k takes only the values 0, 1, 2,

¹⁵Professor Wilf has generously made his book available for free in electronic form. You can download it from <http://www.math.upenn.edu/~wilf/DownldGF.html>.

and 3, with probabilities $p_0, p_1, p_2,$ and $p_3,$ respectively, and no other values. In that case the corresponding generating function would take the form of a cubic polynomial:

$$g(z) = p_0 + p_1z + p_2z^2 + p_3z^3. \quad (12.73)$$

For instance, if we had a network in which nodes of degree 0, 1, 2, and 3 occupied 40%, 30%, 20%, and 10% of the network respectively then

$$g(z) = 0.4 + 0.3z + 0.2z^2 + 0.1z^3. \quad (12.74)$$

As another example, suppose that k follows a Poisson distribution with mean c :

$$p_k = e^{-c} \frac{c^k}{k!}. \quad (12.75)$$

This distribution is generated by the generating function

$$g(z) = \sum_{k=0}^{\infty} e^{-c} \frac{c^k}{k!} z^k = e^{-c} \sum_{k=0}^{\infty} \frac{(cz)^k}{k!} = e^{c(z-1)}. \quad (12.76)$$

Alternatively, suppose that k follows a distribution of the form

$$p_k = Ca^k, \quad (12.77)$$

with $a < 1$, which is an exponential distribution, also sometimes called a geometric distribution in this context. The normalizing constant C is fixed by the requirement that $\sum_k p_k = 1$, which gives $C/(1-a) = 1$ or equivalently $C = 1-a$, and hence

$$p_k = (1-a)a^k. \quad (12.78)$$

Then the generating function is

$$g(z) = (1-a) \sum_{k=0}^{\infty} (az)^k = \frac{1-a}{1-az}, \quad (12.79)$$

so long as $z < 1/a$. If $z \geq 1/a$ the generating function diverges, which seems like a possible concern, but normally we will be interested in generating functions only in the range $0 \leq z \leq 1$ so, given that $a < 1$, the divergence at $1/a$ will not be a problem.

12.10.3 POWER-LAW DISTRIBUTIONS

One special case of particular interest in the study of networks is the power-law distribution. As we saw in Section 10.4, a number of networks, including the

World Wide Web, the Internet, and citation networks, have degree distributions that follow power laws quite closely and this turns out to have interesting consequences that set these networks apart from others. To create and solve models of these networks we would like to be able to write down generating functions for power-law distributions.

Consider, for example, the case of the “pure” power law of Eq. (12.57):

$$p_k = \begin{cases} 0 & \text{for } k = 0, \\ k^{-\alpha} / \zeta(\alpha) & \text{for } k \geq 1, \end{cases} \quad (12.80)$$

where

$$\zeta(\alpha) = \sum_{k=1}^{\infty} k^{-\alpha}, \quad (12.81)$$

is the Riemann zeta function.

For this probability distribution the generating function is

$$g(z) = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} k^{-\alpha} z^k. \quad (12.82)$$

Unfortunately, as discussed in Section 12.8, the sum in this expression cannot be written in closed form, which is somewhat unsatisfactory. As we will see, however, even when expressed as a sum the generating function can still be useful to us.

We should note also that, as mentioned in Section 10.4, degree distributions in real-world networks do not usually follow a power law over their whole range—the distribution is not a pure power law in the sense above. Instead, they typically obey a power law reasonably closely for values of k above some minimum value k_{\min} but below that point have some other behavior. In this case the generating function will take the form

$$g(z) = Q(z) + C \sum_{k=k_{\min}}^{\infty} k^{-\alpha} z^k, \quad (12.83)$$

where $Q(z) = \sum_{k=0}^n p_k z^k$ is a polynomial in z of degree n and C is a normalizing constant. In the calculations in this book we will stick to the pure power-law form, since it illustrates nicely the interesting properties of power-law degree distributions and is relatively simple to work with, but for serious modeling purposes one may sometimes have to use the more complex form of Eq. (12.83).

12.10.4 NORMALIZATION AND MOMENTS

We now look at some of the properties of generating functions that will be useful to us. First of all, note that if we set $z = 1$ in the definition of the

generating function, $g(z) = \sum_k p_k z^k$ (Eq. (12.71)), we get

$$g(1) = \sum_{k=0}^{\infty} p_k. \quad (12.84)$$

Assuming the probability distribution is normalized to unity so that $\sum_k p_k = 1$, this immediately implies that

$$g(1) = 1. \quad (12.85)$$

The derivative of $g(z)$ is

$$g'(z) = \sum_{k=0}^{\infty} k p_k z^{k-1}. \quad (12.86)$$

(We will use the primed notation $g'(z)$ for derivatives of generating functions extensively in this book, as it proves less cumbersome than the more common notation dg/dz .) If we set $z = 1$ in Eq. (12.86) we get

$$g'(1) = \sum_{k=0}^{\infty} k p_k = \langle k \rangle, \quad (12.87)$$

which is the mean value of k . Thus, for example, if p_k is a degree distribution, we can calculate the average degree directly from the generating function by differentiating. This is a very convenient trick. Often we calculate a probability distribution of interest by calculating first its generating function. In principle, we can then extract the distribution itself by applying Eq. (12.72) and so derive any other quantities we want such as averages. But Eq. (12.87) shows us that we don't always have to do this. Some of the quantities we will be interested in can be calculated directly from the generating function without going through any intermediate steps.

In fact, this result generalizes to higher moments of the probability distribution as well. For instance, note that

$$z \frac{d}{dz} \left(z \frac{dg}{dz} \right) = \sum_{k=0}^{\infty} k^2 p_k z^k, \quad (12.88)$$

and, setting $z = 1$, we get

$$\langle k^2 \rangle = \left[\left(z \frac{d}{dz} \right)^2 g(z) \right]_{z=1}. \quad (12.89)$$

It is not hard to show that in general

$$\langle k^m \rangle = \left[\left(z \frac{d}{dz} \right)^m g(z) \right]_{z=1}. \quad (12.90)$$

12.10.5 PRODUCTS OF GENERATING FUNCTIONS

Perhaps the most useful property of generating functions—and the one that makes them important for the study of networks—is the following. Suppose we have m integers k_1, \dots, k_m which are independent random numbers and each is drawn from its own distribution: $p_k^{(1)}$ for k_1 , $p_k^{(2)}$ for k_2 , and so forth. Then the generating function for the probability distribution of the sum $\sum_{i=1}^m k_i$ of these m integers is the *product of the generating functions for the individual distributions*. This is a very powerful result and it is worth taking a moment to see how it arises and what it means.

Given that our integers are independently drawn from their respective distributions, the probability that they take a particular set of values $\{k_i\}$ is simply the product $\prod_i p_{k_i}^{(i)}$ of their individual probabilities, and the probability π_s that the values drawn add up to a specific total s is the sum of this product over all sets $\{k_i\}$ that add up to s :

$$\pi_s = \sum_{k_1=0}^{\infty} \dots \sum_{k_m=0}^{\infty} \delta(s, \sum_i k_i) \prod_{i=1}^m p_{k_i}^{(i)} \tag{12.91}$$

where $\delta(a, b)$ is the Kronecker delta. Then the generating function $h(z)$ for the distribution π_s is

$$h(z) = \sum_{s=0}^{\infty} \pi_s z^s = \sum_{s=0}^{\infty} z^s \sum_{k_1=0}^{\infty} \dots \sum_{k_m=0}^{\infty} \delta(s, \sum_i k_i) \prod_{i=1}^m p_{k_i}^{(i)}. \tag{12.92}$$

Performing the sum over s , we get

$$\begin{aligned} h(z) &= \sum_{k_1=0}^{\infty} \dots \sum_{k_m=0}^{\infty} z^{\sum_i k_i} \prod_{i=1}^m p_{k_i}^{(i)} = \prod_{i=1}^m \sum_{k_i=0}^{\infty} p_{k_i}^{(i)} z^{k_i} \\ &= \prod_{i=1}^m g^{(i)}(z), \end{aligned} \tag{12.93}$$

where

$$g^{(i)}(z) = \sum_{k=0}^{\infty} p_k^{(i)} z^k \tag{12.94}$$

is the generating function for the distribution $p_k^{(i)}$.

Thus the distribution of the sum of a set of independent random integers is generated by the product of their generating functions.

In most of the cases we will study, the random variables of interest to us will all be drawn from the same distribution, say p_k . In this situation all the

generating functions $g^{(i)}(z)$ above are the same function $g(z) = \sum_k p_k z^k$ and Eq. (12.93) simplifies to

$$h(z) = [g(z)]^m. \quad (12.95)$$

That is, the distribution of the sum of m identically distributed random integers is generated by the m th power of the generating function for the individual integers. Thus, for example, if we know the degree distribution of a network, it is a straightforward matter to calculate the probability distribution of the sum of the degrees of m randomly chosen nodes in that network. This will turn out to be important in the developments that follow.

12.10.6 GENERATING FUNCTIONS FOR DEGREE DISTRIBUTIONS

Turning to the application of generating functions in networks, we defined in Section 12.6 the generating functions for the degree distribution and excess degree distribution in the configuration model:

$$g_0(z) = \sum_{k=0}^{\infty} p_k z^k, \quad (12.96)$$

$$g_1(z) = \sum_{k=0}^{\infty} q_k z^k. \quad (12.97)$$

As we pointed out there, these two functions are not really independent, since the excess degree distribution is itself defined in terms of the ordinary degree distribution via Eq. (12.16). We showed that $g_1(z)$ can be calculated from $g_0(z)$ thus:

$$g_1(z) = \frac{g_0'(z)}{g_0'(1)}. \quad (12.98)$$

(See Eq. (12.33).)

As an example, suppose our degree distribution is a Poisson distribution with mean c :

$$p_k = e^{-c} \frac{c^k}{k!}. \quad (12.99)$$

Then its generating function is given by Eq. (12.76) to be

$$g_0(z) = e^{c(z-1)}. \quad (12.100)$$

Applying Eq. (12.98), we then find that

$$g_1(z) = e^{c(z-1)}. \quad (12.101)$$

In other words, $g_0(z)$ and $g_1(z)$ are identical in this case. (This is one reason why calculations are relatively straightforward for the Poisson random graph—there is no difference between the degree distribution and the excess degree distribution, a fact you can easily demonstrate for yourself by substituting Eq. (12.99) directly into Eq. (12.16).)

As another example, a network with degree distribution following the exponential form of Eq. (12.78) has generating functions

$$g_0(z) = \frac{1-a}{1-az}, \quad g_1(z) = \left(\frac{1-a}{1-az}\right)^2, \quad (12.102)$$

while for a power-law distribution, Eq. (12.80), we showed in Section 12.8 that

$$g_0(z) = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} k^{-\alpha} z^k, \quad g_1(z) = \frac{1}{\zeta(\alpha-1)} \sum_{k=1}^{\infty} k^{-(\alpha-1)} z^{k-1}. \quad (12.103)$$

(See Eqs. (12.63) and (12.64).) While these last expressions cannot be written in closed form, they will be useful to us nonetheless.

12.10.7 NUMBER OF SECOND NEIGHBORS OF A NODE

Before we tackle some of the more advanced calculations one can do using generating functions, let us see an example of a calculation we have already performed, redone using generating function methods. We will look at the calculation of the number of second neighbors of a node, which we examined previously in Section 12.5. As we will see, generating functions allow us not only to perform the calculation in a new way, but they also give us new information about the complete distribution of numbers of second neighbors.

We focus on the following question: what is the probability $p_k^{(2)}$ that a node has exactly k second neighbors in a configuration model network? Let us break this probability down by writing it in the form

$$p_k^{(2)} = \sum_{m=0}^{\infty} p_m P^{(2)}(k|m), \quad (12.104)$$

where $P^{(2)}(k|m)$ is the probability of having k second neighbors given that we have m first neighbors and p_m is the ordinary degree distribution, i.e., the distribution of the number of first neighbors. Equation (12.104) says that the total probability of having k second neighbors is the probability of having k second neighbors given that we have m first neighbors, averaged over all possible values of m . We assume that we are given the degree distribution p_m . Our goal is to find $P^{(2)}(k|m)$ and then complete the sum.

This turns out to be a hard calculation to do directly, so instead of calculating the probability $p_k^{(2)}$ itself, we calculate its generating function:

$$g^{(2)}(z) = \sum_{k=0}^{\infty} p_k^{(2)} z^k = \sum_{k=0}^{\infty} \sum_{m=0}^{\infty} p_m P^{(2)}(k|m) z^k = \sum_{m=0}^{\infty} p_m \left[\sum_{k=0}^{\infty} P^{(2)}(k|m) z^k \right]. \quad (12.105)$$

Note that the quantity in brackets on the right is itself the generating function for $P^{(2)}(k|m)$. This generating function turns out to be easy to calculate.

As illustrated in Fig. 12.2 on page 383, the number of second neighbors of a node is equal to the sum of the excess degrees of the first neighbors, and the excess degrees are distributed according to the excess degree distribution q_k , Eq. (12.16). Now we can use the product property of generating functions derived in Section 12.10.5: the sum of a set of m random integers that all have the same distribution is generated by the m th power of the generating function for each individual one. Thus the generating function for $P^{(2)}(k|m)$ is simply the m th power of the generating function $g_1(z)$ for the excess degree distribution:

$$\sum_{k=0}^{\infty} P^{(2)}(k|m) z^k = [g_1(z)]^m. \quad (12.106)$$

Substituting this result back into Eq. (12.105), we now have

$$g^{(2)}(z) = \sum_{m=0}^{\infty} p_m [g_1(z)]^m. \quad (12.107)$$

But now we notice an interesting fact: the sum in this expression is nothing other than the generating function $g_0(z)$ for the degree distribution, Eq. (12.96), evaluated at the point $g_1(z)$. In other words,

$$g^{(2)}(z) = g_0(g_1(z)). \quad (12.108)$$

So once we know the generating functions for our two basic degree distributions, the generating function for the distribution of second neighbors is simple to calculate.

Armed with this result, we can now calculate various things about the second neighbors. For instance, we can calculate the average number of second neighbors using Eq. (12.87), which says that the average of a probability distribution is given by the first derivative of its generating function evaluated at $z = 1$. The derivative of $g^{(2)}(z)$ is

$$\frac{dg^{(2)}}{dz} = g_0'(g_1(z)) g_1'(z). \quad (12.109)$$

Setting $z = 1$ and recalling that $g_1(1) = 1$ (Eq. (12.85)), we find that the average number c_2 of second neighbors is

$$c_2 = g'_0(1)g'_1(1). \quad (12.110)$$

But $g'_0(1)$ is itself a first derivative evaluated at $z = 1$, and hence is equal to the average of the distribution p_k , i.e., to the average degree of the network $\langle k \rangle$, a result that you can easily confirm for yourself from the definition (12.96). Similarly, $g'_1(1)$ is the average excess degree. We calculated the average excess degree previously in Eq. (12.20), but let us do the calculation again explicitly using the generating function:

$$\begin{aligned} g'_1(1) &= \sum_{k=0}^{\infty} kq_k \\ &= \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} k(k+1)p_{k+1} = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} (k-1)kp_k \\ &= \frac{1}{\langle k \rangle} (\langle k^2 \rangle - \langle k \rangle), \end{aligned} \quad (12.111)$$

where we have used the definition of the excess degree distribution from Eq. (12.16).

Putting these results together, the mean number c_2 of second neighbors of a node can be written

$$c_2 = \langle k \rangle \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle} = \langle k^2 \rangle - \langle k \rangle, \quad (12.112)$$

which is identical to our previous result, Eq. (12.19).

While this calculation merely duplicates an earlier result, the generating function method can also tell us new things. For instance, we can now calculate the precise probability that a node has exactly k second neighbors by differentiating $g^{(2)}(z)$ according to Eq. (12.72). As an example, the probability $p_0^{(2)}$ of having no second neighbors at all is given by

$$p_0^{(2)} = g^{(2)}(0) = g_0(g_1(0)) = g_0(q_0) = g_0(p_1/\langle k \rangle), \quad (12.113)$$

where we have used the definitions (12.16) and (12.97) of the excess degree distribution and its generating function $g_1(z)$. On a network with a Poisson degree distribution, for instance (see Eq. (12.99) and following), this gives

$$p_0^{(2)} = e^{c(e^{-c}-1)}. \quad (12.114)$$

This is a result that would not have been easy to calculate using traditional methods.

We can take these calculations further: the same method can be used to calculate the probability distribution of the number of third neighbors, or indeed neighbors at any distance d . The interested reader can find these and other results in Ref. [369], but for the moment we will move on to other things.

12.10.8 GENERATING FUNCTIONS FOR THE SMALL COMPONENTS

A good example of the power of the generating function method is in calculating the properties of the small components in the configuration model. We will go through this calculation in some detail, showing how the method can be used to calculate the entire distribution of sizes of the small components. The calculation involves some significantly more complex mathematics than we have seen so far.

The fundamental quantity we will focus on is the probability π_s that a randomly chosen node belongs to a small (non-giant) component of size s . We will calculate this probability by first calculating its generating function

$$h_0(z) = \sum_{s=1}^{\infty} \pi_s z^s. \quad (12.115)$$

Note that the minimum value of s is 1, since every node belongs to a component of size at least one (namely the node itself).

Consider again Fig. 12.4a on page 392, which shows the neighborhood of a node i in one of the small components. As before, let us imagine removing node i from the network along with all of its edges, as shown in Fig. 12.4b. If we do this, the shaded areas in the figure become separate components in their own right, and the size of the complete component to which node i belongs in the original network is equal to the sum of the sizes of these new components, plus one for node i itself.

A crucial point to note, however, is that the neighbors n_1, n_2, \dots of node i are, by definition, reached by following an edge. Hence, as we have discussed, these are not typical network nodes, being more likely to have high degree than the typical node. Thus, the sizes of the components they belong to in Fig. 12.4b—the shaded regions in the figure—are not in general distributed according to π_s . Instead they have some other distribution. Let us denote this distribution by ρ_s . More specifically, let ρ_s be the probability that the node at the end of an edge belongs to a small component of size s after that edge is

removed. Let us also define the generating function for this distribution thus:

$$h_1(z) = \sum_{s=1}^{\infty} \rho_s z^s. \quad (12.116)$$

We don't yet know the value of ρ_s or its generating function and we will have to calculate them later, but for the moment let us proceed with the information we have.

Suppose that node i on the original network has degree k and let us denote by $P(s|k)$ the probability that, after i is removed, its k neighbors belong to small components of sizes summing to exactly s . Alternatively, $P(s-1|k)$ is the probability that i itself belongs to a small component of size s given that its degree is k . Then the total probability π_s that i belongs to a small component of size s is this probability averaged over k thus:

$$\pi_s = \sum_{k=0}^{\infty} p_k P(s-1|k), \quad (12.117)$$

where p_k is the degree distribution, as usual.¹⁶ Substituting this expression into Eq. (12.115) we then get an expression for the generating function for π_s as follows:

$$\begin{aligned} h_0(z) &= \sum_{s=1}^{\infty} \sum_{k=0}^{\infty} p_k P(s-1|k) z^s = z \sum_{k=0}^{\infty} p_k \sum_{s=1}^{\infty} P(s-1|k) z^{s-1} \\ &= z \sum_{k=0}^{\infty} p_k \sum_{s=0}^{\infty} P(s|k) z^s. \end{aligned} \quad (12.118)$$

The final sum $\sum_s P(s|k)z^s$ in this expression is the generating function for the probability that, after we have removed node i from the network, its k neighbors belong to small components whose size sums to s . But now we can use the product property of generating functions derived in Section 12.10.5, which tells us that the generating function for this sum is just equal to the k th power of the

¹⁶We here average over the degree distribution p_k for the whole network. In Section 12.7.1 we pointed out that the nodes in the small components have a degree distribution different from that of the network as a whole, and one might imagine we should use the same modified degree distribution in Eq. (12.117) too. However, the quantity π_s in Eq. (12.117) is a probability for the entire network, not just the nodes in the small components—it is the probability that *any* node, including one in the giant component, belongs to a small component of size s . Obviously a node in the giant component doesn't belong to a small component of any size, and this is factored into the calculation as we will see—the sum of the probabilities π_s over all s is not equal to 1 for exactly this reason.

generating function for the size of the component any single neighbor belongs to—the function that we denoted $h_1(z)$ in Eq. (12.116). Thus,

$$h_0(z) = z \sum_{k=0}^{\infty} p_k [h_1(z)]^k = z g_0(h_1(z)). \quad (12.119)$$

We still don't know the generating function $h_1(z)$ but we can derive it now quite easily. We consider again the network in which node i is removed and ask what the value is of the probability ρ_s that one of the neighbors of i belongs to a component of size s in this network. In the limit of large network size, the removal of the single node i will have no effect on the degree distribution, so the network still has the same distribution as before, which means that if a neighbor has degree k then its probability of belonging to a component of size s is $P(s-1|k)$, just as before. Note, however, that since the neighbor was reached by following an edge from i , its degree, discounting the edge to i that has been removed, follows the excess degree distribution q_k of Eq. (12.16). So the analog of Eq. (12.117) is now

$$\rho_s = \sum_{k=0}^{\infty} q_k P(s-1|k), \quad (12.120)$$

and, substituting this expression into Eq. (12.116), we have

$$h_1(z) = \sum_{s=1}^{\infty} \sum_{k=0}^{\infty} q_k P(s-1|k) z^s = z \sum_{k=0}^{\infty} q_k \sum_{s=0}^{\infty} P(s|k) z^s. \quad (12.121)$$

As before, the last sum is the generating function for $P(s|k)$, which is equal to $[h_1(z)]^k$, and hence

$$h_1(z) = z \sum_{k=0}^{\infty} q_k [h_1(z)]^k = z g_1(h_1(z)). \quad (12.122)$$

Collecting together our results, the generating functions for π_s and ρ_s thus satisfy

$$h_0(z) = z g_0(h_1(z)), \quad (12.123)$$

$$h_1(z) = z g_1(h_1(z)). \quad (12.124)$$

If we can solve the second of these equations for $h_1(z)$ then we can substitute the result into the first and we have our answer for $h_0(z)$ —the generating function for the complete probability distribution of the sizes of small components.

12.10.9 COMPLETE DISTRIBUTION OF SMALL COMPONENT SIZES

In practice, it is often not easy to solve Eq. (12.124) for $h_1(z)$, and even if it is, extracting the actual component size distribution from the generating function using Eq. (12.72) could still be difficult. It transpires, however, that we don't need to do these things to find the size distribution. In a surprising and beautiful turn of events, it turns out we can calculate the complete distribution of component sizes even in cases where the generating functions themselves cannot be calculated, using methods drawn from the calculus of complex variables.

The first thing to note is that, since a component cannot have size zero, the generating function for the probabilities π_s has the form

$$h_0(z) = \sum_{s=1}^{\infty} \pi_s z^s, \quad (12.125)$$

with the sum starting at 1. Dividing by z and differentiating $s - 1$ times, we then find that

$$\pi_s = \frac{1}{(s-1)!} \left[\frac{d^{s-1}}{dz^{s-1}} \left(\frac{h_0(z)}{z} \right) \right]_{z=0} \quad (12.126)$$

(which is just a minor variation on the standard formula, Eq. (12.72)). Using Eq. (12.123), this can also be written

$$\begin{aligned} \pi_s &= \frac{1}{(s-1)!} \left[\frac{d^{s-1}}{dz^{s-1}} g_0(h_1(z)) \right]_{z=0} \\ &= \frac{1}{(s-1)!} \left[\frac{d^{s-2}}{dz^{s-2}} [g'_0(h_1(z))h'_1(z)] \right]_{z=0}. \end{aligned} \quad (12.127)$$

Now we make use of the Cauchy formula for the n th derivative of a function $f(z)$ of a complex variable, which says that

$$\frac{d^n f}{dz^n} \Big|_{z=z_0} = \frac{n!}{2\pi i} \oint \frac{f(z)}{(z-z_0)^{n+1}} dz, \quad (12.128)$$

where the integral goes counterclockwise around a contour in the complex plane that encloses z_0 but no poles in $f(z)$. Applying this formula to Eq. (12.127) with $z_0 = 0$ we get

$$\pi_s = \frac{1}{2\pi i(s-1)} \oint \frac{g'_0(h_1(z))}{z^{s-1}} \frac{dh_1}{dz} dz. \quad (12.129)$$

For our contour, we choose an infinitesimal circle around the origin.

Changing the integration variable to h_1 , we can also write this as

$$\pi_s = \frac{1}{2\pi i(s-1)} \oint \frac{g'_0(h_1)}{z^{s-1}} dh_1, \quad (12.130)$$

where we are regarding z now as a function of h_1 , rather than the other way around. It is important that the contour followed by h_1 surrounds the origin, so let us pause for a moment to demonstrate that it does. Expanding Eq. (12.124) around $z = 0$, we find that

$$h_1(z) = z g_1(h_1(0)) + O(z^2) = z g_1(0) + O(z^2), \quad (12.131)$$

where we have made use of the fact that $h_1(0) = 0$, since components cannot have size zero—see Eq. (12.116). In the limit of small $|z|$ where the terms of order z^2 can be neglected, Eq. (12.131) implies that h_1 traces an infinitesimal circle about the origin if z does, since the two are proportional to one another. Moreover, they rotate the same way about the origin, since the constant of proportionality $g_1(0) = q_0$ is positive, so the sign of our integral is correct.¹⁷

Now we make use of Eq. (12.124) to eliminate z from Eq. (12.130) and write

$$\begin{aligned} \pi_s &= \frac{1}{2\pi i(s-1)} \oint \frac{[g_1(h_1)]^{s-1} g'_0(h_1)}{h_1^{s-1}} dh_1 \\ &= \frac{g'_0(1)}{2\pi i(s-1)} \oint \frac{[g_1(h_1)]^s}{h_1^{s-1}} dh_1, \end{aligned} \quad (12.132)$$

where we have made use of Eq. (12.33) in the second line. Given that the contour surrounds the origin, this integral is now in the form of Eq. (12.128) again, and hence

$$\pi_s = \frac{\langle k \rangle}{(s-1)!} \left[\frac{d^{s-2}}{dz^{s-2}} [g_1(z)]^s \right]_{z=0}, \quad (12.133)$$

where we have written $g'_0(1) = \langle k \rangle$.

The only exception to this formula is for the case $s = 1$, for which Eq. (12.129) gives $0/0$ and is therefore clearly incorrect. However, since the only way to belong to a component of size 1 is to have no connections to any other nodes, the probability π_1 is trivially equal to the probability of having degree zero:

$$\pi_1 = p_0. \quad (12.134)$$

¹⁷One might worry that q_0 could be zero, in which case $h_1(z) = O(z^2)$ and h_1 would go *twice* around the origin when z goes around once. This possibility we can rule out, however. Given that $q_0 = p_1/\langle k \rangle$ (see Eq. (12.16)), $q_0 = 0$ would imply that the network has no nodes of degree 1. But the small components are, as we have said, trees, and all trees must have at least one node of degree 1 (unless they consist of only a single node of degree zero). A simple way to see this is to note that the mean degree of a tree is strictly less than 2 (because the mean degree is $2m/n = 2(n-1)/n < 2$). Thus, if there are any small components at all of size greater than 1, we must have $q_0 > 0$.

Equations (12.133) and (12.134) give the probability that a randomly chosen node belongs to a component of size s in terms of the degree distribution. In principle if we know p_k we can now calculate π_s . It is not always easy to perform the derivatives in practice and sometimes we may not even know the generating function $g_1(z)$ in closed form, but at least in some cases the calculations are possible.

As an example, consider a network with the exponential (or geometric) degree distribution

$$p_k = (1 - a) a^k, \quad (12.135)$$

with $a < 1$. The generating functions $g_0(z)$ and $g_1(z)$ are then given by Eq. (12.102) and it is not hard to show that

$$\frac{d^n}{dz^n} [g_1(z)]^s = \frac{(2s - 1 + n)!}{(2s - 1)!} \frac{[g_1(z)]^s}{(a^{-1} - z)^n}, \quad (12.136)$$

and hence that

$$\pi_s = \frac{(3s - 3)!}{(s - 1)!(2s - 1)!} a^{s-1} (1 - a)^{2s-1}. \quad (12.137)$$

Figure 12.8 shows a comparison of this formula with the results of numerical simulations for the case $a = 0.3$ and, as we can see, the agreement between formula and simulations is excellent (even though the simulated network is necessarily finite in size while the calculations are performed in the limit of large n).

Generating functions provide a powerful tool for the calculation of a wide range of network properties. We will see a number of further examples in this book.

12.11 OTHER RANDOM GRAPH MODELS

In this chapter and the previous one we have studied in detail the two most fundamental network models, the Poisson random graph and the configuration model. There are, however, many other random graph models in addition to these, including models of directed networks, bipartite networks, acyclic networks, networks with degree correlations, clustering, and community structure, and many others. In the remainder of this chapter we look briefly at a selection of these models.

12.11.1 DIRECTED NETWORKS

As discussed in Section 6.4, many networks, including the World Wide Web, metabolic networks, food webs, and others, are directed. The Poisson random

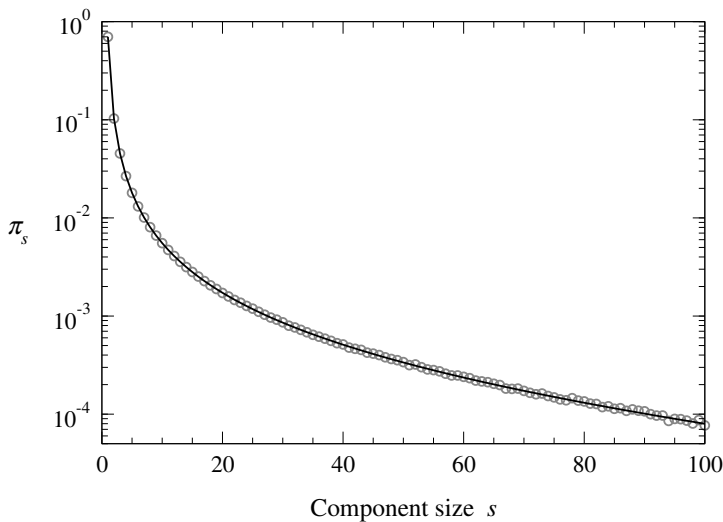


Figure 12.8: The distribution of component sizes in a configuration model. The probability π_s that a node belongs to a component of size s for the configuration model with an exponential degree distribution of the form (12.135) for $a = 0.3$. The solid lines represent the exact formula, Eq. (12.137), for the $n \rightarrow \infty$ limit and the points are measurements of π_s averaged over 100 computer-generated networks with $n = 10^7$ nodes each.

graph $G(n, p)$ can be generalized to directed networks in a straightforward fashion: independently with probability p one places between each pair of nodes i, j a directed edge from i to j and, with the same probability p , another directed edge from j to i —see Exercise 11.7 on page 367. More useful in practice is the directed generalization of the configuration model, which is also straightforward: we specify an in- and out-degree for every node, place the corresponding number of ingoing and outgoing stubs at each node, then repeatedly choose pairs of stubs at random, one ingoing and one outgoing, and join them to create a directed edge, until no unused stubs remain. The result is a matching of the stubs drawn uniformly at random from the set of all possible matchings, just as in the undirected configuration model. The only small catch is that we must make sure that the total numbers of ingoing and outgoing stubs in the network are the same, so that no stubs are left over at the end of the matching process.

We can calculate many of the same properties for the directed version of the configuration model as for the undirected version, such as the probability

of an edge between two nodes, for example. The probability that a particular outgoing stub at node j attaches to one of the k_i^{in} ingoing stubs at node i is

$$\frac{k_i^{\text{in}}}{\sum_j k_j^{\text{in}}} = \frac{k_i^{\text{in}}}{m}, \quad (12.138)$$

where m is the total number of edges and we have made use of Eq. (6.19). Since the total number of outgoing stubs at j is k_j^{out} , the total expected number of directed edges from node j to node i is then

$$p_{ij} = \frac{k_i^{\text{in}} k_j^{\text{out}}}{m}, \quad (12.139)$$

which is also the probability of an edge from i to j in the limit of a large sparse network. This is similar to the corresponding result, Eq. (12.2), for the undirected configuration model, but not identical—note that there is no factor of two now in the denominator.

As in the undirected case we can, if we prefer, work with the degree distribution, rather than the degree sequence, by first drawing a set of degrees from a specified degree distribution and then using those degrees as the starting point for the matching process above. Again, one must ensure that the total numbers of ingoing and outgoing stubs are the same, meaning that the sum of in-degrees over the whole network must equal the sum of out-degrees.

As discussed in Section 10.3, the most correct way to describe the degree distribution of a directed network is as a joint distribution: we define p_{jk} to be the fraction of nodes in the network that have in-degree j and out-degree k . This allows for the possibility that the in- and out-degrees of nodes are correlated. For instance, it would allow us to represent a network in which the in- and out-degrees of each node were exactly equal to one another. (This is a rather extreme form of correlation, but it demonstrates the point.)

By methods analogous to those for the undirected case, it is also possible to calculate properties of the directed model such as the average number of neighbors a certain distance away from a node, whether there is a giant component in the network, and the size of the giant component if there is one. Note that, as described in Section 6.12.1, we distinguish in directed networks between strongly and weakly connected components, and both can be studied using the directed version of the configuration model. The generating function methods of Section 12.10 also extend to the directed model, though the generating functions are now functions of two variables, not just one. For details, see Refs. [149,369].

Given the joint degree distribution we can still, if we wish, calculate the distributions of in- or out-degrees alone, which are given by $\sum_k p_{jk}$ and $\sum_j p_{jk}$ respectively.

12.11.2 BIPARTITE NETWORKS

It is similarly straightforward to generalize the configuration model to the case of bipartite networks—networks with two types of nodes and edges only between unlike types (see Section 6.6). We create two sets of nodes and assign a degree to each node, represented by stubs in the usual way. Then we repeatedly pick one stub at random from each of the two sets and join the stubs together to form a complete edge. When all stubs have been used up the network is complete. This again generates a matching of stubs drawn uniformly from the set of all possible matchings. A requirement of the model is that the total number of stubs attached to nodes of each type be the same, so that no stubs get left over at the end of the matching process. This means that the sums of the degrees of the nodes of each type must be equal.

Again we can calculate many of the same quantities for this model as for the ordinary configuration model. There are now two degree distributions for the network, one for each of the types of nodes, and two excess degree distributions. Otherwise, however, calculations proceed along similar lines to those presented in this chapter. One can, for instance, calculate edge probabilities and component sizes, as well as some properties unique to bipartite networks, such as properties of the one-mode projections onto nodes of a single type (Section 6.6.1). See Ref. [369] for details.

12.11.3 ACYCLIC NETWORKS

A directed network is acyclic if it contains no closed directed loops. As shown in Section 6.4.1, an equivalent statement is that the nodes of the network can be arranged in a line such that all edges go in one direction only along the line. Acyclic networks occur most commonly in situations where there is a natural such one-dimensional ordering that dictates the directions of edges, such as time ordering in the case of a citation network: all citations between scientific papers go from later papers to earlier ones, so if the papers are arranged in time order all edges will point the same way—backward in time.

To create an acyclic version of the configuration model, one starts by arranging n nodes in some order, which we can think of as time ordering if we wish. We assign an in- and out-degree to each node just as in the directed model of Section 12.11.1, and place the corresponding number of ingoing and outgoing stubs at each node. Then we work through the nodes one after another, in order from “earliest” to “latest,” and connect each outgoing stub we encounter to an ingoing stub chosen uniformly at random from among the currently unused ingoing stubs at earlier nodes. When we have gone through all nodes in this

manner the network is complete.

One important point to note about this model is that only certain degree sequences will give rise to legitimate networks. In order for the process above to work, the out-degree of the i th node in the sequence cannot be greater than the number of unused stubs at all earlier nodes; otherwise some outgoing stubs would have nowhere to attach to. In mathematical terms, this means that we require $k_i^{\text{out}} \leq \sum_{j=1}^{i-1} k_j^{\text{in}} - \sum_{j=1}^{i-1} k_j^{\text{out}}$, or equivalently

$$\sum_{j=1}^{i-1} k_j^{\text{in}} \geq \sum_{j=1}^i k_j^{\text{out}}, \quad (12.140)$$

for all i . Not only is this a necessary condition for the process to work, it is also clearly a sufficient one—if it is true then the process is always possible. Add to this the one further condition that the sum of all in-degrees must equal the sum of all out-degrees (as is true for all directed networks), then we have a complete set of criteria for deciding when a degree sequence is allowed and when it isn't.

Random directed acyclic networks of this kind have been studied to some extent [255], but not that many of their properties are known—calculations for these networks appear to be harder in many ways than for other types of random graphs.

12.11.4 DEGREE CORRELATIONS

A crucial property of most real-world networks, discussed in Section 7.7.3, is assortativity by degree, the tendency of nodes to connect to others with degree either similar to or different from their own. There is no simple generative process, equivalent to the stub-matching procedures of previous sections, to create random networks with correlated degrees. However, one can still calculate many properties of such networks analytically using generalizations of the methods of this chapter. The trick is to specify a joint excess degree distribution $q_{kk'}$, the probability that two nodes connected by an edge have excess degrees k and k' respectively. Then the ordinary excess degree distribution q_k is given by

$$q_k = \sum_{k'=0}^{\infty} q_{kk'} \quad (12.141)$$

and the probability that a neighbor of node i has excess degree k' , given that i itself has excess degree k , is $q_{kk'}/q_k$. Armed with these quantities we can calculate a range of properties of correlated networks, including average

numbers of neighbors, clustering coefficients, and component properties. See Refs. [211,350,457] for details.

12.11.5 CLUSTERING AND TRANSITIVITY

As discussed in Sections 7.3 and 10.6, many observed networks, especially social networks, show high levels of transitivity—two nodes are more likely to be connected if they share a common neighbor—and networks often contain a larger-than-expected number of triangles as a result. This is an easy property to incorporate into our random graph models: we simply create random triangles in the network in much the same way that we previously created random edges. In the simplest version of this idea, one assigns two separate degree-like parameters to each node, one representing the number of single edges attached to the node and the other representing the number of corners of triangles attached to the node. One then assigns edge stubs and corner stubs to nodes in the appropriate numbers and performs two different matching processes. In the first process one picks pairs of edge stubs at random and joins them to form edges in the usual manner. In the second process one picks trios of corner stubs and creates triangles between the three nodes they attach to—see Fig. 12.9. Many properties of this model can be calculated exactly in the limit of large n by methods analogous to those for the configuration model—see Ref. [360]. The model can also be extended in a straightforward manner to include not only triangles but other motifs as well, such as groups of four or five nodes, connected in various ways [256].

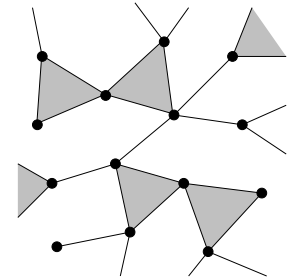


Figure 12.9: A random graph model of a network with clustering. In this model one separately places single edges between pairs of nodes and complete triangles between trios of nodes.

12.11.6 ASSORTATIVE MIXING AND COMMUNITY STRUCTURE

In Sections 7.7 and 10.7 we discussed the phenomenon of assortative mixing, observed particularly in social networks, whereby nodes that share some feature or attribute in common are more likely to be connected. People of the same age, income, nationality, race, or educational level, for example, are typically more likely to be friends than people who differ on these things. In rarer cases one also sees *disassortative* mixing, where nodes are more likely to be connected if they don't have the same features.

The standard random graph model of this kind of assortative mixing is the *stochastic block model*, which is an assortatively mixed version of the Poisson random graph of Chapter 11 (not the configuration model). In this model one takes n nodes and divides them into some number q of groups or types, numbered from 1 to q , which might represent languages, age brackets, ethnicities, educational levels, or any other variables of interest. Then one places

undirected edges independently at random between node pairs, just as in the ordinary random graph, except that the probabilities of the edges now depend on the groups nodes belong to. Specifically, instead of all edges having the same probability, we define a set of quantities p_{rs} which give the probability of an edge between two nodes, one of which is in group r and the other in group s .

The quantities p_{rs} form a $q \times q$ matrix whose diagonal elements p_{rr} represent the probability of edges within groups and off-diagonal elements represent probabilities between different groups. The matrix is symmetric since the probability p_{rs} of an edge between groups r and s is by definition the same as the probability p_{sr} between groups s and r .

If the diagonal elements of the matrix are larger than the off-diagonal ones, then edges will be more likely inside groups than between them and we get a network with traditional assortative mixing. The model can generate disassortative mixing too, however, if we make the diagonal elements smaller than the off-diagonal ones, although it is less often used this way.

Though it is simple (and well studied), the stochastic block model suffers from many of the same shortcomings as the Poisson random graph (see Section 11.8), particularly that within any group the degrees of the nodes have a Poisson distribution, which is very different from the degree distributions observed in real-world networks. To rectify this problem one can create a generalization of the block model analogous to the configuration model. As in the configuration model one chooses the degree k_i of each node, which is represented by stubs attached to the node, but one also chooses the number m_{rs} of edges that will fall between groups r and s (or within groups in the case where $r = s$). Then one goes through the edges in turn and, for each edge that is required to fall between groups r and s , we pick at random one stub from group r and one from group s and join them to form an edge. We repeat this process until all edges have been created.

Note that for this model to work, the numbers of edges m_{rs} must match the chosen degree sequence. Specifically, the sum of the degrees of nodes in group r must equal the number of ends of edges that are required to attach to group r .

This model has the nice features of being easy to describe and a true extension of the configuration model to the case of assortative mixing. However, in practice it is rarely used. Instead, most mathematical work on this problem has used a different model, the *degree-corrected stochastic block model*, which is an extension not of the configuration model but of the model of Chung and Lu described in Section 12.1.2, in which we fix not the degrees of the nodes in the network but the expected degrees.

Recall that in the Chung–Lu model we place edges independently between node pairs with probabilities $c_i c_j / 2m$, where c_i is the expected degree of node i —the actual degree can vary around this value, only the mean being fixed at c_i . In the corresponding block model, we place edges independently with probabilities $\omega_{rs} c_i c_j / 2m$, where r and s are the groups to which nodes i and j belong. Thus ω_{rs} plays the role of a factor that modifies the probability of an edge, making it bigger or smaller for different groups, relative to its value in the Chung–Lu model. If the diagonal elements ω_{rr} of the (symmetric) matrix formed by these parameters are bigger than the off-diagonal ones, then again we have traditional assortative mixing in the network; if they are smaller we have disassortative mixing.

Note that we do not have complete freedom about how the parameters ω_{rs} are chosen if we want c_i to be equal to the expected degree of node i . The expected degree of node i is just the sum of the expected numbers of edges from that node to all nodes. If we denote by g_i the group to which node i belongs, then the expected degree is equal to

$$\sum_j \omega_{g_i g_j} \frac{c_i c_j}{2m} = \frac{c_i}{2m} \sum_j \omega_{g_i g_j} c_j. \quad (12.142)$$

If we want this to be equal to c_i for all nodes i , then we must have

$$\sum_j \omega_{r g_j} c_j = 2m \quad (12.143)$$

for all r . This places q separate linear constraints on the values of the quantities ω_{rs} , one for each of the q groups.

The models described in this section can be used as models of networks with assortative mixing, but their primary use is in the different but related domain of “community detection.” In Chapter 14 we will look in some detail at the phenomenon of community structure in networks—the commonly observed division of networks into groups or communities of nodes such that there are dense connections within groups but only sparser connections between groups. This is similar to the phenomenon of assortative mixing discussed above except that, in most cases, there is no known external variable or characteristic (such as age or income or nationality) dictating the groups—they are simply an observed feature of the network. Nonetheless, such groups can be a useful guide to the structure and function of many networks: the breakdown of the network into groups can give us insight into how it formed or the dynamics of interactions between its nodes, or it can just provide a useful way to cluster the network for visualization purposes.

For these and other reasons, it is often useful to be able to pick out the groups or communities in a network, and a range of different “community detection” algorithms have been proposed that aim to do exactly this. One of the most elegant approaches uses the stochastic block model or its degree-corrected variant: we fit the model to the data for the observed network and the parameters of the fit tell us about the community structure in the network in much the same way that the fit of a straight line through a set of data points can tell us about the slope of the data. Community detection methods, including this one, are described at length in Chapter 14.

12.11.7 DYNAMIC NETWORKS

Another variant on the random graph idea incorporates dynamics that change the structure of the network over time. As pointed out in Section 6.7, most real-world networks are not static but evolve over time. Social networks change as people make new friends or fall out of touch with old ones; the Web changes when pages or links are added or deleted; biological networks, such as metabolic networks, change over evolutionary time. What is the equivalent of models like the Poisson random graph and the configuration model for this dynamic case? There is no one answer to this question. Many models for dynamic networks have been proposed—see [239] for a survey. However, perhaps the simplest and most direct equivalent of the Poisson random graph is the *Markov model* in which edges appear and disappear between nodes uniformly and independently at random [217, 483]. In this model there is some probability λ per unit time that an edge will appear between a pair of nodes where currently there is none, and another probability μ per unit time that an existing edge will disappear. Over time, as edges come and go, the network changes shape, but at any particular moment in time it takes the form of a Poisson random graph, since each edge is equally likely to exist at any time. We can calculate the random graph probability p that two nodes are connected by an edge at any particular time by noting that the mean time before an edge appears in a currently vacant spot is $1/\lambda$, while the mean time for which an edge exists before disappearing is $1/\mu$. Hence the average fraction of the time for which a specific edge exists, which is also the probability p of an edge, is

$$p = \frac{1/\mu}{1/\lambda + 1/\mu} = \frac{\lambda}{\lambda + \mu}. \quad (12.144)$$

Normally, we are interested in sparse networks for which $p \ll 1$, meaning that $\lambda \ll \mu$ and edges do not last very long. Given that the average timescale on which an edge disappears is $1/\mu$, this is also the average timescale on

which the whole network turns over: wait this long and most edges will have disappeared, being replaced by others.

One can make a dynamic version of the configuration model as well, or more accurately a dynamic version of the model of Chung and Lu [103] (see Section 12.1.2). Since this model can contain multiedges we use a slightly different dynamic process. With probability λ per unit time a new edge appears between a given pair of nodes whether or not there is an edge (or several) already there, and with probability μ per unit time existing edges vanish. The equilibrium number p_{ij} of edges between nodes i and j can then be calculated by equating the average rates of appearance and disappearance of edges between those nodes, meaning that $\lambda = \mu p_{ij}$, so that $p_{ij} = \lambda/\mu$. If we now choose the value of λ separately for each pair of nodes to be $\lambda_{ij} = \mu c_i c_j / 2m$, where c_i and c_j are the desired degrees of the nodes, then at any moment the expected number of edges between i and j is

$$p_{ij} = \frac{\lambda_{ij}}{\mu} = \frac{c_i c_j}{2m}, \quad (12.145)$$

which is the correct value for the Chung–Lu model. This ensures that the degrees of the nodes will indeed have average values c_i (although they, like everything else, fluctuate over time), and at any particular moment the network has the form of the Chung–Lu model.

This still leaves one free parameter μ , which, as before, controls the rate of turnover of edges in the network. If we wait for time $1/\mu$, the network will still be a Chung–Lu style network, with the same expected degrees at the nodes, but the structure of the network itself will have been erased and replaced with a new one.

Relatively few studies have been made of the mathematics of networks like these, although in principle they have some interesting properties. For instance, the component structure of the network changes over time along with everything else, so even if two nodes are not in the same component at the present moment, they may be at some later time. One could study, therefore, how long we have to wait before two nodes that are currently unconnected will become connected. This could have implications, for instance, for the spread of disease over contact networks.

12.11.8 THE SMALL-WORLD MODEL

A model with a somewhat different style and purpose is the *small-world model* proposed by Watts and Strogatz [466]. The small-world model is a stylized model originally intended to illustrate how two characteristic features of social

networks—high clustering coefficient and short path lengths—can coexist in the same network. Today, arguably, we have better models of networks with these features, such as the models in Section 12.11.5, and in hindsight the small-world model's main contribution may be a different one: to show why the small-world effect (the existence of short paths between most node pairs) is so prevalent in networks of all kinds (see Sections 4.6 and 10.2).

The model is defined as follows. One starts off with a regular lattice of some kind. Versions of the model have been explored that make use of various lattices, but the original model used a one-dimensional lattice as depicted in Fig. 12.10a. The nodes are arranged in a line and each node is connected by edges to the c nodes nearest to it, where for consistency c should be an even number. To make analytic treatment easier, we typically apply periodic boundary conditions to the line, effectively bending it around into a circle, as in Fig. 12.10b.

Now we take this network and randomize it by moving or *rewiring* some of the edges from their positions around the circle to new random positions. Specifically, we go through each of the edges around the circle in turn and with some probability p we remove that edge and replace it with one that joins two nodes chosen uniformly at random from the network.¹⁸ The result is shown in Fig. 12.10c. The randomly placed edges are commonly referred to as *shortcuts* because, as shown in the figure, they create shortcuts from one part of the circle to another.

If all the edges in the circle were rewired randomly in this way, then the result would be a standard Poisson random graph of the kind we studied in Chapter 11: all edges are placed uniformly at random. Thus, the parameter p in the small-world model interpolates between the circular lattice we started out with and the random graph. When $p = 0$ no edges are rewired and we retain the original circle. When $p = 1$ all edges are rewired and we have a random graph. The original purpose of the model was to argue that while the random graph has short path lengths but low clustering coefficient and the circular lattice has high clustering but long paths, there is a significant parameter range in between in which the network has both properties simultaneously, short paths and high clustering.

¹⁸In fact, in the original small-world model, as defined by Watts and Strogatz [466], only one end of each edge—say the more clockwise end—was rewired and the other left where it was. Some other constraints were also imposed, such as the constraint that no two edges connect the same node pair. These choices, however, make little difference to the behavior of the model in practice, but do make calculations more difficult, so in most studies the simpler version of the model defined here has been used.

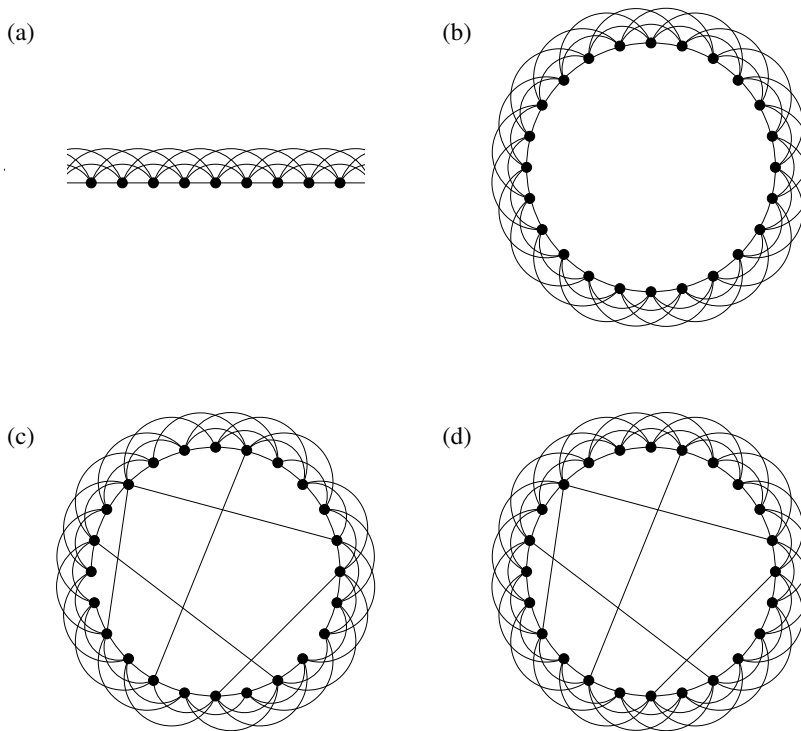


Figure 12.10: Construction of the small-world model. (a) Nodes are arranged on a line and each is connected to its c nearest neighbors, where $c = 6$ in this example. (b) Applying periodic boundary conditions makes the line into a circle. (c) Independently with probability p edges are randomly “rewired,” meaning that they are removed from the circle and placed between randomly chosen nodes, creating shortcuts across the circle. In this example $n = 24$, $c = 6$, and $p = 0.07$, so that 5 out of 72 edges are rewired in this fashion. (d) In the alternate version of the model only the shortcuts are added and no edges are removed from the circle.

But there is also another lesson, perhaps more salient, that we can draw from the small-world model. It turns out that short path lengths develop in the model even for very small values of the probability p . Only a tiny fraction of the edges need be rewired for the path lengths to become short. If we define the small-world effect as we did in Section 11.7 to mean that the diameter of the network increases no faster than $\log n$ as $n \rightarrow \infty$, then it turns out that this effect is achieved in the small-world model when only a vanishing fraction of all edges have been rewired [370].

This offers some explanation of why the small-world effect is so ubiquitous

in networks of all kinds. It tells us that if one randomizes even a vanishing fraction of the edges in a network then the small-world effect appears, which in turn implies that almost all networks show the small-world effect. If they did not—if, say, only a half of networks showed the small-world effect—then when we randomized edges we would fail to get the small-world effect about half the time. Since this does not happen, we conclude that the small-world effect is almost universal.

Unfortunately, it is hard to demonstrate this result rigorously using the small-world model as defined above because the model is difficult to treat by analytic means. For this reason the model is often studied in a slightly different variant which is easier to treat [370]. In this variant, shown in Fig. 12.10d, edges are added between randomly chosen node pairs just as before, but no edges are removed from the original circle. This leaves the circle intact, which makes calculations much simpler. For ease of comparison with the original small-world model, the definition of the parameter p is kept the same: for every edge in the circle we add with independent probability p an additional shortcut between two nodes chosen uniformly at random.¹⁹

A downside of this version of the model is that it no longer becomes a random graph in the limit $p = 1$. Instead it becomes a random graph plus the original circle. This, however, turns out not to be a significant problem, since most of the interest in the model lies in the regime where p is small and in this regime the two models differ hardly at all; the only difference is the presence in the second variant of a small number of edges around the circle that would be absent in the first, having been rewired.

EXERCISES

12.1 As described in Section 12.1, the configuration model can be thought of as the ensemble of all possible matchings of edge stubs, where node i has k_i stubs. Show that for a given degree sequence the number of matchings is $(2m)!/(2^m m!)$, which is independent of the degree sequence, except for the dependence on the total number of edges m .

¹⁹Equivalently, one could just say that the number of shortcuts added is drawn from a Poisson distribution with mean $\frac{1}{2}ncp$.

12.2 As discussed in Section 12.1, the configuration model generates each possible matching of edge stubs with equal probability. It does not, however, generate each possible network with equal probability, because different networks correspond to different numbers of matchings. One can generate all the matchings that correspond to a given network by taking any one matching for that network and permuting the stubs at each node in every possible way. Since the number of permutations of the k_i stubs at a node i is $k_i!$, this seems to imply that the number of matchings corresponding to each network is $N = \prod_i k_i!$, which takes the same value for all networks, since the degrees are fixed. However, this is not completely correct. If a network contains self-edges or multiedges then not all permutations of the stubs in the network result in a new matching. Show that the actual number of matchings corresponding to a given network is

$$N = \frac{\prod_i k_i!}{\prod_{i < j} A_{ij}! \prod_i A_{ii}!!},$$

where $n!!$ is the so-called double factorial of n , which is equal to $n(n-2)(n-4)\dots 2$ when n is even and $0!! = 1$. If there are no multiedges or self-edges, the expression above reduces to the earlier formula $\prod_i k_i!$ but in the general case it does not.

12.3 The “friendship paradox” of Section 12.2 says that your friends tend to have more friends than you do on account of the fact that they are reached by following an edge and hence have higher-than-average expected degree. The *generalized friendship paradox* is a related phenomenon in which your friends are richer/smarter/happier than you are (or any other characteristic). This happens when characteristics are correlated with degree. If, for instance, rich people tend to have more friends, then wealth and degree will be positively correlated in the friendship network. Since your friends have higher degree than you on average, we can then expect them also to be richer.

- a) Suppose we have some value, such as wealth, on each node in a configuration model network. Let us denote the value on node i by x_i . Show that if we follow any edge in the network, the average value of x_i on the node at its end will be

$$\langle x \rangle_{\text{edge}} = \frac{1}{2m} \sum_i k_i x_i.$$

- b) Hence show that the equivalent of Eq. (12.14)—the difference between the average value of x_i for a friend and the average value for the network as a whole—is

$$\langle x \rangle_{\text{edge}} - \langle x \rangle = \frac{\text{cov}(k, x)}{\langle k \rangle},$$

where $\text{cov}(k, x) = \langle kx \rangle - \langle k \rangle \langle x \rangle$ is the covariance of k and x over nodes.

Equation (12.14) is the special case of this result where x_i is the degree of the node, so that $\text{cov}(k, x) = \sigma_k^2$, the variance of the degree distribution.

12.4 Consider a configuration model in which every node has the same degree k .

- a) What is the degree distribution p_k ? What are the generating functions g_0 and g_1 for the degree distribution and the excess degree distribution?

- b) Show that the giant component fills the whole network for all $k \geq 3$.
- c) What happens when $k = 1$?
- d) The case $k = 2$ is significantly harder to analyze. If you are feeling ambitious, try showing that in the limit of large n the probability π_s that a node belongs to a component of size s is given by $\pi_s = 1/[2\sqrt{n(n-s)}]$.

12.5 Show that in a configuration model network with nodes of degree 2 and greater, but no nodes of degree 0 or 1, there are no small components (or, more properly, the fraction of nodes belonging to such components tends to zero as $n \rightarrow \infty$).

12.6 Write a computer program in the programming language of your choice that generates a configuration model network with nodes of degree 1 and 3 only and then calculates the size of the largest component.

- a) Use your program to calculate the size of largest component for a network of $n = 10\,000$ nodes with $p_1 = 0.6$ and $p_3 = 0.4$ (and $p_k = 0$ for all other values of k).
- b) Modify your program to calculate the size of the largest component for values of p_1 from 0 to 1 in steps of 0.01, then make a graph of the results as a function of p_1 . Hence estimate the value of p_1 at the phase transition where the giant component disappears. Compare your result to the predictions of the analytic calculation in Section 12.6.1.

12.7 Consider the binomial probability distribution $p_k = \binom{n}{k} p^k (1-p)^{n-k}$.

- a) Show that the probability generating function for this distribution is $g(z) = (pz + 1 - p)^n$.
- b) Find the first and second moments of the distribution from Eqs. (12.87) and (12.89) and hence show that the variance of the distribution is $\sigma^2 = np(1-p)$.
- c) Show that the sum k of two numbers drawn independently from the same binomial distribution is distributed according to $\binom{2n}{k} p^k (1-p)^{2n-k}$.

12.8 Generating functions can be used to study many things in addition to probability distributions.

- a) The Fibonacci numbers $1, 1, 2, 3, 5, 8, \dots$ have the definitive property that each is the sum of the previous two. The generating function for the Fibonacci numbers is the power series whose coefficients are the Fibonacci numbers: $f(z) = 1 + z + 2z^2 + 3z^3 + 5z^4 + \dots$. Show that $f(z) = 1/(1 - z - z^2)$. What is the largest (most positive) z for which the series defining the generating function converges to a finite value?
- b) A sequence of numbers a_k with $k = 1, 2, 3, \dots$ satisfies the recurrence relation

$$a_k = \begin{cases} 1 & \text{for } k = 1, \\ \sum_{j=1}^{k-1} a_j a_{k-j} & \text{for } k > 1. \end{cases}$$

Show that the generating function $g(z) = \sum_{k=1}^{\infty} a_k z^k$ satisfies $g(z) = \frac{1}{2}(1 - \sqrt{1 - 4z})$.

12.9 Starting from the generating function $h_0(z)$ in Eq. (12.123), or otherwise, do the following.

- a) Derive an expression in terms of g_0 and g_1 for the mean-square size of the component in a configuration model network to which a randomly chosen node belongs;
 b) Show that for a Poisson degree distribution with mean degree c , the mean-square size is $1/(1-c)^3$ in the regime where there is no giant component.

12.10 Consider a configuration model network with degree distribution $p_k = 2^{-(k+1)}$ for all $k \geq 0$, in the limit of large n .

- a) Show that this degree distribution is correctly normalized.
 b) Derive expressions for the generating functions $g_0(z)$ and $g_1(z)$ for the degree distribution and the excess degree distribution.
 c) Calculate the average degree c_1 of a node and the average number of second neighbors c_2 .
 d) Does the network have a giant component? How do you know?
 e) What is the probability that a randomly chosen node belongs to a component of size 3?

12.11 Consider the configuration model with degree distribution $p_k = (1-a)a^k$ with $a < 1$, so that the generating functions $g_0(z)$ and $g_1(z)$ are given by Eq. (12.102).

- a) Show that the probability u of Eq. (12.30) satisfies the cubic equation

$$a^2 u^3 - 2a u^2 + u - (1-a)^2 = 0.$$

- b) Noting that $u = 1$ is always a trivial solution of this equation, show that the non-trivial solution corresponding to the existence of a giant component satisfies the quadratic equation $a^2 u^2 - a(2-a)u + (1-a)^2 = 0$.
 c) Hence show that as a fraction of the size of the network, the size of the giant component, if there is one, is

$$S = \frac{3}{2} - \sqrt{a^{-1} - \frac{3}{4}}.$$

- d) Show that the giant component exists only if $a > \frac{1}{3}$.

12.12 Consider a configuration model with degree distribution p_k in the limit of large n .

- a) Write an expression for the probability that a node of degree k belongs to the giant component, in terms of the quantity u of Eq. (12.30).
 b) Hence derive an expression for the probability that a node in the giant component has degree k .
 c) Show that the average degree of a node in the giant component is $c(1-u^2)/S$, where c and S are, as usual, the average degree in the network as a whole and the size of the giant component.
 d) Show that $\sum_{jk} p_j p_k (j-k)(u^j - u^k) \leq 0$, and hence that $\sum_k k u^k p_k \leq \sum_k k p_k \sum_k u^k p_k$. By rewriting the latter inequality in terms of the generating functions g_0 and g_1 of Eqs. (12.26) and (12.29), prove that in any configuration model network the probability that a node of degree 2 belongs to the giant component is greater than or equal to the probability that the average node belongs to the giant component.

- e) By combining the results from (c) and (d), show that the average degree of nodes in the giant component of any configuration model network is greater than or equal to the average degree in the network as a whole.

12.13 Using the results of Section 12.10.9, show that in a random graph with a Poisson degree distribution the probability π_s that a randomly chosen node belongs to a component of size s is $\pi_s = (cs)^{s-1}e^{-cs}/s!$ where c is the mean degree.

12.14 Consider a network model in which edges are placed independently between each pair of nodes i, j with probability $p_{ij} = Kf_i f_j$, where K is a constant and f_i is a number assigned to node i . Show that the expected degree c_i of node i in this model is proportional to f_i and hence that the only possible choice of edge probability with this form is $p_{ij} = c_i c_j / 2m$, as in the Chung–Lu model of Section 12.1.2.

12.15 Consider the example model discussed in Section 12.6.1, a configuration model with nodes of degree three or less only and generating functions given by Eqs. (12.34) and (12.35).

- a) In the regime in which there is no giant component, show that the average size of the component to which a randomly chosen node belongs is

$$\langle s \rangle = 1 + \frac{(p_1 + 2p_2 + 3p_3)^2}{p_1 - 3p_3}.$$

- b) In the same regime find the probability that a randomly chosen node belongs to components of size 1, 2, and 3.

12.16 Consider a configuration model with degree distribution $p_k = Cka^k$, where a and C are positive constants and $a < 1$.

- Calculate the value of the constant C as a function of a .
- Calculate the mean degree of the network.
- Calculate the mean-square degree of the network.
- Hence, or otherwise, find the value of a that marks the phase transition between the region in which the network has a giant component and the region in which it does not. Does the giant component exist for larger or smaller values than this?

You may find the following sums useful in performing the calculations:

$$\sum_{k=0}^{\infty} ka^k = \frac{a}{(1-a)^2}, \quad \sum_{k=0}^{\infty} k^2 a^k = \frac{a+a^2}{(1-a)^3}, \quad \sum_{k=0}^{\infty} k^3 a^k = \frac{a+4a^2+a^3}{(1-a)^4}.$$

12.17 Suppose the Internet is found to have a power-law degree distribution $p_k \sim k^{-\alpha}$ for $k \geq 1$, with $\alpha \simeq 2.5$.

- Make a mathematical model of the Internet using the configuration model with this degree distribution. Write down the fundamental generating functions g_0 and g_1 . (The generating functions cannot be written in closed form, so you should leave them in sum form.)

- b) Hence estimate what fraction of the nodes on the Internet you expect to be functional at any one time (where functional means they can actually send data over the network to most other nodes).

12.18 Consider a directed random graph of the kind discussed in Section 12.11.1, where the degree distribution is specified by fixing the joint probability p_{jk} that a node has in-degree j and out-degree k . By analogy with the argument we made for the configuration model, a node in such a network has a giant out-component if the number of nodes reachable from it, by following edges in their forward direction, grows exponentially as we get further and further from the node.

- a) By calculating the average rate of growth (or decay) of the number of reachable nodes with distance, show that if the in- and out-degrees of nodes are uncorrelated, i.e., if p_{jk} factors into a product of separate probabilities for j and k , a giant out-component exists if and only if $c > 1$, where c is the mean degree of the network, either in or out. This is the equivalent of the criterion in Eq. (12.24) for the ordinary configuration model.
- b) Show that the same condition also applies for the existence of a giant in-component.
- c) Argue that the existence of both a giant in-component and a giant out-component at the same time also implies the existence of giant weakly and strongly connected components, and hence that all the giant components—in, out, strong, and weak—have the same condition for their existence.
- d) In real directed networks the degrees are usually correlated (or anti-correlated). The correlation can be quantified by the covariance ρ of in- and out-degrees. Show that in the presence of correlations, the condition for the existence of the giant components is $c(c - 1) + \rho > 0$.
- e) In the World Wide Web the in- and out-degrees of nodes have a measured covariance of about $\rho = 180$. The mean degree is around $c = 4.6$. On the basis of these numbers, do we expect the Web to have giant components?

12.19 Consider a bipartite analog of the configuration model, as described in Section 12.11.2, in which there are two types of nodes, A and B, and edges run only between nodes of unlike types. Each node type can have its own degree distribution—there is no need for the distributions to be the same or even similar.

- a) Depending on the exact form of the degree distributions, the network may or may not contain a giant component. Derive a condition in terms of the mean and mean-square degrees of the two types, equivalent to Eq. (12.24) for the ordinary configuration model, that tells us when a giant component exists.
- b) Define u_A to be the probability that the node of type A at the end of an edge is *not* in the giant component, and similarly for u_B and nodes of type B. Show that $u_A = g_1^A(u_B)$ and $u_B = g_1^B(u_A)$ where g_1^A and g_1^B are the generating functions for the excess degrees of nodes of type A and B respectively.
- c) Give an expression for the fraction S_A of nodes of type A in the giant component.

CHAPTER 13

MODELS OF NETWORK FORMATION

A discussion of models of the formation of networks, particularly networks that grow by addition of nodes, such as the World Wide Web or citation networks

THE MODELS described in Chapters 11 and 12 provide good tools for studying the structural features of networks, such as components, degree distributions, path lengths, and so forth. Moreover, as we will see in later chapters, they can also serve as a convenient starting point for further modeling work, such as the modeling of network resilience or the spread of disease.

But there is another important class of network models that has an entirely different purpose. In the models we have seen so far, the parameters of the network, such as the numbers of nodes and edges or the degree distribution, are fixed from the outset, chosen by the modeler to have some desired values. For instance, if we are interested in networks with power-law degree distributions, we can make a random graph model with a power-law degree distribution as in Section 12.8 and then explore its structure analytically or computationally. But models of this kind offer no explanation of *why* the network should have a power-law degree distribution in the first place. In this chapter we describe models of a different kind that do offer such an explanation.

The models in this chapter are models of the mechanisms by which networks are created. The idea behind these models is to explore hypothesized mechanisms of network formation to see what structures they produce. If the structures are similar to those of networks we observe in the real world, it suggests—though does not prove—that similar mechanisms may be at work in

the real networks.

The best-known example of such a network model, and the one that we study first in this chapter, is the “preferential attachment” model for the growth of networks with power-law degree distributions. Later in the chapter we examine a number of other models, including node copying models and network optimization models.

13.1 PREFERENTIAL ATTACHMENT

As discussed in Section 10.4, many networks are observed to have degree distributions that approximately follow power laws, at least in the tail of the distribution. Examples include the Internet, the World Wide Web, citation networks, and some social and biological networks. The power law is a somewhat unusual distribution and its occurrence in empirical data is often considered a potential indicator of interesting underlying processes. A natural question to ask therefore is how might a network come to have such a distribution? This question was first directly considered in the 1970s by Price [394], who proposed a simple and elegant model of network formation that gives rise to power-law degree distributions.

Price was interested in, among other things, the citation networks of scientific papers, having authored an important early paper on the topic in the 1960s in which he pointed out the power-law degree distribution seen in these networks [393]. In considering the possible origins of the power law, Price was inspired by the work of economist Herbert Simon [429], who noted the occurrence of power laws in a variety of (non-network) economic data, such as the distribution of people’s personal wealth. Simon proposed an explanation for the wealth distribution based on the idea that people who already have a lot of money gain more at a rate proportional to how much they currently have. This seems a reasonable supposition. Wealthy individuals make money primarily by investing their wealth, and the return on their investment is essentially proportional to the amount invested. Simon was able to show mathematically that this “rich-get-richer” effect can give rise to a power-law distribution and Price adapted Simon’s methods, with relatively little change, to the network context. Price gave a name to Simon’s mechanism, calling it *cumulative advantage*,¹ although today it is more often known as *preferential attachment*, a name coined later by Barabási and Albert [40]. In this book we principally use the latter term.

¹Simon himself called the mechanism the *Yule process*, in recognition of the statistician Udny Yule, who had studied a simple version many years earlier [478].

For a discussion of the general properties of power laws see Refs. [335] and [357].

See Section 3.2 for a discussion of citation networks.

Price's model of a citation network is as follows. Papers are published continually (though they do not have to be published at a constant rate) and newly appearing papers cite previously existing ones. As discussed in Section 3.2, the papers and citations form a directed citation network, the papers being the nodes and the citations being the directed edges between them. Since no paper ever disappears after it is published, nodes in this network are created but never destroyed.

Let the mean number of papers cited by a newly appearing paper be c . In the language of networks, c is the average out-degree of the citation network. In the language of publishing it is the average size of the bibliography of a paper. The model allows the actual sizes of bibliographies to fluctuate around this average from paper to paper. So long as the distribution of sizes satisfies a few basic sanity conditions,² only the average value is important for the behavior of the model in the limit of large network size. In real citation networks the sizes of bibliographies also vary from one field to another and depend on when papers were published, the average bibliography having grown larger over the years in most fields, but these effects are neglected in the model.

The crucial central assumption of Price's model is that the papers cited by a newly appearing paper are chosen at random *with probability proportional to the number of citations they already have*. Choosing papers at random is clearly not an accurate representation of what happens in the real publication process. Price's model ignores such important issues as which papers are most relevant topically or which papers are most original or best written or the difference between research articles and reviews, or any of the many other factors that certainly affect real citation patterns. The model is thus very much a simplified portrait of the citation process. As we have seen with the random graphs of previous chapters, however, even simple models can lead to real insights. While we certainly need to remember that the model represents only one aspect of the citation process—and a hypothetical one at that—let us press on and see what we can discover.

As with personal wealth, it is not implausible that the number of citations a paper receives could increase with the number it already has. When one reads papers, one often looks up the other works that those papers cite and reads some of them too. If a work is cited often, then, all other things being equal, we are more likely to come across it than a less cited work. And if we read it and

²The main condition on the distribution is that it should have finite variance. This rules out, for example, cases in which bibliographies have a power-law distribution of sizes with exponent less than 3. Empirical evidence suggests that real bibliographies have an unexceptionable distribution of sizes with a modest and finite variance, so the assumptions of Price's model are met.

like it, then perhaps we will cite it ourselves if we write a paper on the same topic. This does not mean that the probability of a paper receiving a citation is precisely proportional to the number of citations the paper has already, but it does at least give some justification for why the rich should get richer in this paper citation context.

In fact, upon further reflection, it's clear that the probability of receiving a new citation cannot be *precisely* proportional to the number of citations a paper already has. Except in unusual circumstances, papers start out life with zero citations, which, with a strict proportionality rule, would mean that their probability of getting new citations would also be zero and so they would have zero citations for ever afterwards. To get around this hitch, Price proposed that the probability that a paper receives a new citation should be proportional to the number that it already has plus a positive constant a . (In fact, Price only considered one special case $a = 1$ in his original paper, but there seems to be no particular reason to limit ourselves to this case, so we will treat the case of general $a > 0$.)

The constant a in effect gives each paper a number of “free” citations to get it started in the race—each paper acts as though it started off with a citations instead of zero. This gives all papers a chance to accrue new citations, even if they currently have none. An alternative interpretation is that a certain fraction of citations go to papers chosen uniformly at random without regard for how many citations they currently have, while the rest go to papers chosen in proportion to current citation count. (We discuss this interpretation in more detail in Section 13.1.2, where we use it to construct a fast algorithm for simulating Price's model.)

We also need to specify the starting state of the network, the nucleus from which it grows. It turns out, in fact, that in the limit of large network size the predictions of the model don't depend on the starting state, but we could, for instance, start the network with a small set of initial papers having zero citations each.

Thus, in summary, Price's model consists of a growing network of papers and their citations in which nodes (papers) are continually added but none are ever taken away, each paper cites on average c others (so that the mean out-degree is c), and the cited papers are chosen at random³ with probability

³There is nothing in the definition of Price's model to prevent a paper from listing the same other paper twice in its bibliography, something that doesn't happen in real citation networks. Such double citations would correspond to directed multiedges in the citation network (see Section 6.1) while true citation networks are simple networks having no multiedges. However, as with the configuration model of Chapter 12, the probability of generating a multiedge vanishes in the limit

proportional to their in-degree plus a constant a .

One important property of Price's model is immediately apparent: it generates purely acyclic networks (see Section 6.4.1), since every edge points from a more recently added node to a less recently added one, i.e., backward in time. Thus all directed paths in the network point backward in time and hence there can be no closed loops, because to close a loop we would need edges pointing forward in time as well. This fits well with the original goal of the model as a model of citation, since citation networks are acyclic, or very nearly so (see Section 3.2). On the other hand, it fits poorly with some other directed networks such as the World Wide Web, although the model is still sometimes used as a model for power-law distributions in the Web.

13.1.1 DEGREE DISTRIBUTION OF PRICE'S MODEL

Armed with our definition of Price's model, we will now write down equations governing the distribution of the in-degrees of nodes, i.e., the numbers of citations received by papers in terms of the parameters c and a , and hence solve for the degree distribution, at least in the limit of large network size. We will discuss models of both directed and undirected networks in this chapter, so we will need to be careful to distinguish in-degree in the directed case from ordinary undirected degree in the undirected case. Previously in this book we have done this by denoting the in-degree of a node i by k_i^{in} (see Section 6.10), but this notation can make our equations quite difficult to read, so in the interests of clarity we will in this chapter adopt instead the notation introduced by Dorogovtsev *et al.* [148] in which the in-degree of node i is denoted q_i . Degrees in undirected networks will still be denoted k_i just as before.

So consider Price's model of a growing network and let $p_q(n)$ be the fraction of nodes in the network that have in-degree q when the network contains n nodes—this is the in-degree distribution of the network—and let us examine what happens when we add a single new node to the network.

Consider one of the citations made by this new node. Following the definition of the model, the probability that the citation is to a particular other node i is proportional to $q_i + a$, where a is a positive constant. Since the citation in question has to be to *some* node, this probability must be normalized such that its sum over all i is 1. In other words, the correctly normalized probability must

of large network size, so the predictions of the model in this limit are not altered by allowing them, and doing so makes the mathematical treatment of the model much simpler.

be

$$\frac{q_i + a}{\sum_i (q_i + a)} = \frac{q_i + a}{n\langle q \rangle + na} = \frac{q_i + a}{n(c + a)}, \quad (13.1)$$

where we have written the average in-degree as $\langle q \rangle = n^{-1} \sum_i q_i$. In the second equality we have made use of the fact that the average in-degree is equal to the average out-degree in any directed network (see Eq. (6.20)), and that the average out-degree of this network is c by definition, so $\langle q \rangle = c$.

Each newly appearing paper cites c others on average, so the expected number of new citations to node i upon appearance of our new paper is c times Eq. (13.1). And there are $np_q(n)$ nodes with in-degree q in our network and hence the expected number of new citations to all nodes with in-degree q is

$$np_q(n) \times c \times \frac{q + a}{n(c + a)} = \frac{c(q + a)}{c + a} p_q(n). \quad (13.2)$$

Now we can write down a so-called *master equation* for the evolution of the in-degree distribution as follows. When we add a single new node to our network of n nodes, the number of nodes in the network with in-degree q increases by one for every node previously of in-degree $q - 1$ that receives a new citation,⁴ thereby becoming a node of in-degree q . From Eq. (13.2) we know that the expected number of such nodes is

$$\frac{c(q - 1 + a)}{c + a} p_{q-1}(n). \quad (13.3)$$

Similarly, we lose one node of in-degree q every time such a node receives a new citation, thereby becoming a node of in-degree $q + 1$. The expected number of times this happens is

$$\frac{c(q + a)}{c + a} p_q(n). \quad (13.4)$$

The number of nodes with in-degree q in the network after the addition of a single new node is $(n + 1)p_q(n + 1)$ which, putting together the results above, is given by

$$(n + 1)p_q(n + 1) = np_q(n) + \frac{c(q - 1 + a)}{c + a} p_{q-1}(n) - \frac{c(q + a)}{c + a} p_q(n). \quad (13.5)$$

The first term on the right-hand side here represents the number of nodes previously of in-degree q , the second term represents the nodes gained, and the third term represents the nodes lost.

⁴In theory, it also increases by one if a node of in-degree $q - 2$ receives two new citations, and similarly for larger numbers of citations. This, however, would create a multiedge, and multiedges, as we have said, are vanishingly improbable in the limit of large network size, so we can ignore this possibility.

Equation (13.5) applies for all values of q except $q = 0$. When $q = 0$ there are no nodes of lower degree that can gain an edge to become nodes of degree zero, so the second term in Eq. (13.5) doesn't appear. On the other hand, we gain a node of degree zero whenever a new node is added to the network, since papers have no citations when they are first published. Since exactly one node is added in going from a network of n nodes to a network of $n + 1$, the appropriate equation for $q = 0$ is

$$(n + 1)p_0(n + 1) = np_0(n) + 1 - \frac{ca}{c + a}p_0(n). \quad (13.6)$$

Now let us consider the limit of large network size $n \rightarrow \infty$ and calculate the asymptotic form of the degree distribution in this limit.⁵ Taking the limit $n \rightarrow \infty$ and using the shorthand $p_q = p_q(\infty)$, Eqs. (13.5) and (13.6) become

$$p_q = \frac{c}{c + a} [(q - 1 + a)p_{q-1} - (q + a)p_q] \quad \text{for } q \geq 1, \quad (13.7)$$

$$p_0 = 1 - \frac{ca}{c + a}p_0 \quad \text{for } q = 0. \quad (13.8)$$

The second of these equations we can easily rearrange to give an explicit expression for the fraction p_0 of degree-zero nodes:

$$p_0 = \frac{1 + a/c}{a + 1 + a/c}. \quad (13.9)$$

The solution for $q \geq 1$ is a little more complicated, though only a little. Rearranging Eq. (13.7) for p_q we find that

$$p_q = \frac{q + a - 1}{q + a + 1 + a/c} p_{q-1}. \quad (13.10)$$

We can use this equation to calculate p_q iteratively for all values of q starting from our solution for p_0 , Eq. (13.9). First, we set $q = 1$ in Eq. (13.10) to get

$$p_1 = \frac{a}{a + 2 + a/c} p_0 = \frac{a}{(a + 2 + a/c)} \frac{(1 + a/c)}{(a + 1 + a/c)}. \quad (13.11)$$

Now we can use this result to calculate p_2 :

$$p_2 = \frac{a + 1}{a + 3 + a/c} p_1 = \frac{(a + 1)a}{(a + 3 + a/c)(a + 2 + a/c)} \frac{(1 + a/c)}{(a + 1 + a/c)}, \quad (13.12)$$

⁵Strictly, we should first prove that the degree distribution *has* an asymptotic form for large n and doesn't go on changing forever, but for the purposes of the present discussion let us assume that there is an asymptotic form.

and

$$p_3 = \frac{a+2}{a+4+a/c} p_2 = \frac{(a+2)(a+1)a}{(a+4+a/c)(a+3+a/c)(a+2+a/c)} \frac{(1+a/c)}{(a+1+a/c)}, \quad (13.13)$$

and so forth. It's easy to see that for general q the correct expression is

$$p_q = \frac{(q+a-1)(q+a-2)\dots a}{(q+a+1+a/c)\dots(a+2+a/c)} \frac{(1+a/c)}{(a+1+a/c)}. \quad (13.14)$$

This is effectively a complete solution for the degree distribution of Price's model, but there is a little more we can do to write it in a useful form. We make use of the gamma function,

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt, \quad (13.15)$$

which has the useful property that⁶

$$\Gamma(x+1) = x\Gamma(x) \quad (13.16)$$

for all $x > 0$. Iterating this formula, we see that

$$\frac{\Gamma(x+n)}{\Gamma(x)} = (x+n-1)(x+n-2)\dots x. \quad (13.17)$$

Using this result in Eq. (13.14) we can write

$$p_q = (1+a/c) \frac{\Gamma(q+a)\Gamma(a+1+a/c)}{\Gamma(a)\Gamma(q+a+2+a/c)}. \quad (13.18)$$

This expression can be simplified further by writing it in terms of Euler's beta function, which is defined by

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}. \quad (13.19)$$

⁶This result can be proved using integration by parts:

$$\Gamma(x+1) = \int_0^\infty t^x e^{-t} dt = -[t^x e^{-t}]_0^\infty + x \int_0^\infty t^{x-1} e^{-t} dt = x\Gamma(x),$$

where the boundary term [...] disappears at both limits.

If we multiply both the numerator and the denominator of Eq. (13.18) by $\Gamma(2 + a/c) = (1 + a/c)\Gamma(1 + a/c)$, we find that

$$p_q = \frac{\Gamma(q + a)\Gamma(2 + a/c)}{\Gamma(q + a + 2 + a/c)} \times \frac{\Gamma(a + 1 + a/c)}{\Gamma(a)\Gamma(1 + a/c)}, \quad (13.20)$$

or

$$p_q = \frac{B(q + a, 2 + a/c)}{B(a, 1 + a/c)}. \quad (13.21)$$

Note that this expression is not only correct for $q \geq 1$ but also gives the correct value when $q = 0$.

Equation (13.21) is sometimes known as the *Yule distribution*, following the work of Udny Yule [478] in the 1920s, who derived it by different methods. One of the nice things about this form is that it depends on q only via the first argument of the upper beta function. Thus, if we want to understand the shape of the degree distribution we only need to understand the behavior of this one function. In particular, let us examine the behavior for large q and fixed a and c . For large values of its first argument, we can rewrite the beta function using Stirling's approximation for the gamma function [2]

$$\Gamma(x) \simeq \sqrt{2\pi} e^{-x} x^{x-\frac{1}{2}}, \quad (13.22)$$

which means that

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)} \simeq \frac{e^{-x} x^{x-\frac{1}{2}}}{e^{-(x+y)} (x+y)^{x+y-\frac{1}{2}}} \Gamma(y). \quad (13.23)$$

But

$$(x+y)^{x+y-\frac{1}{2}} = x^{x+y-\frac{1}{2}} \left(1 + \frac{y}{x}\right)^{x+y-\frac{1}{2}} \simeq x^{x+y-\frac{1}{2}} e^y, \quad (13.24)$$

where the last equality becomes exact in the limit of large x . Then

$$B(x, y) \simeq \frac{e^{-x} x^{x-\frac{1}{2}}}{e^{-(x+y)} x^{x+y-\frac{1}{2}} e^y} \Gamma(y) = x^{-y} \Gamma(y). \quad (13.25)$$

In other words, the beta function $B(x, y)$ falls off as a power law for large values of x , with exponent y .

Applying this finding to Eq. (13.21) we then discover that for large values of q the degree distribution of our network goes as $p_q \sim (q + a)^{-\alpha}$, or simply

$$p_q \sim q^{-\alpha} \quad (13.26)$$

when $q \gg a$, where the exponent α is

$$\alpha = 2 + \frac{a}{c}. \quad (13.27)$$

Thus Price's model for a citation network gives rise to a degree distribution with a power-law tail. This is very much in keeping with the degree distributions of real citation networks, which, as we saw in Fig. 10.8c on page 323, appear to have clear power-law tails.

Note that the exponent $\alpha = 2 + a/c$ is strictly greater than 2 (since a and c are both strictly positive). Most measurements put the exponent of the power law for citation networks around $\alpha = 3$ (see Table 10.1), which is easily achieved in the model by setting the constants a and c equal. In a typical experimental situation the exponent α and the parameter c , the mean size of a paper's bibliography, are easily measured, but the parameter a , which represents the number of "free" effective citations a paper receives upon publication, is not. Typically therefore the value of a is extracted by rearranging Eq. (13.27) to give $a = c(\alpha - 2)$.

While it is intriguing that Price's simple model generates a power-law degree distribution similar to that seen in real networks, we should not take the details of the model too seriously, nor the exact relationship between the parameters and the exponent of the power law. As we noted at the start of this section, the model is highly simplified and substantially incomplete as a model of the citation process, omitting many factors that are undoubtedly important for real citations, including the quality and relevance of papers, developments and fashions in the field of study, the reputation of the publishing journal and of the author, and many others besides. Still, Price's model is striking in its ability to reproduce one of the most interesting features of citation networks using only a small number of reasonable assumptions, and many scholars believe that it may capture the fundamental mechanism behind the observed power-law degree distribution.

13.1.2 COMPUTER SIMULATION OF PRICE'S MODEL

When Price proposed his model in 1976, analytic treatments like the one in the previous section were essentially the only tool available for understanding the behavior of such models. Today, however, we can go further and study the operation of the model explicitly by performing computer simulations following the rules Price laid out. In addition to providing a useful check on our solution for the degree distribution, such simulations also allow us to generate real examples of networks on our computer. We can then measure these networks to determine the values, within the model, of any network quantities we like—path lengths, correlations, clustering coefficients, and so forth—including ones for which we do not at present have an analytic solution. Researchers have also made use of simulated networks as a convenient but still

relatively realistic substrate for other kinds of calculations, including solutions of dynamical models, percolation processes, opinion formation models, and others.

The random graph models of Chapters 11 and 12 are simple to implement on a computer. A direct translation of the models' definition into computer code will generate networks in an accurate and efficient manner. For Price's model, however, an efficient computer implementation requires a little more work. At first sight, the problem appears straightforward. Typically, one simulates the model with the out-degrees of all nodes fixed to be exactly equal to c , where c is restricted to integer values. (In the original model and in the analysis of the previous section, c was only the average out-degree—actual out-degree could fluctuate about the average and need not be an integer.) Then the only complicated part of the simulation is the selection of the nodes that receive new edges, which must be done in a random but non-uniform way as a function of the nodes' current in-degree. There are standard techniques for simulating such non-uniform random processes and one can without too much labor create a program that carries out the steps of the model.

This, however, is not usually the best way to proceed. A naive direct simulation of this kind becomes slow when the network gets large, because of the way the random node selection process works. This limits the size of the networks that can be generated. Luckily, there is a much faster way to perform the simulation that allows large networks to be generated quickly, while still being simple to program on a computer. This method, first proposed by Krapivsky and Redner [279], works as follows.

When we create a new edge in Price's model we attach it to a node chosen in proportion to in-degree plus a constant a . Let us denote by θ_i the probability that an edge attaches to node i , which from Eq. (13.1) is given by

$$\theta_i = \frac{q_i + a}{n(c + a)}. \quad (13.28)$$

Now consider an alternative process in which upon creating a new edge we do one of two things. With some probability ϕ we attach the edge to a node chosen strictly in proportion to its current in-degree, i.e., with probability

$$\frac{q_i}{\sum_j q_j} = \frac{q_i}{nc}. \quad (13.29)$$

Alternatively, with probability $1 - \phi$, we attach to a node chosen uniformly at random from all n possibilities, i.e., with probability $1/n$. Then the total probability θ'_i of attaching to node i in this process is

$$\theta'_i = \phi \frac{q_i}{nc} + (1 - \phi) \frac{1}{n}. \quad (13.30)$$

Now let us make the choice $\phi = c/(c + a)$, so that

$$\theta'_i = \frac{c}{c + a} \frac{q_i}{nc} + \left(1 - \frac{c}{c + a}\right) \frac{1}{n} = \frac{q_i + a}{n(c + a)}. \quad (13.31)$$

This, however, is precisely equal to the probability θ_i , Eq. (13.28), of selecting a node in the Price model and the two processes thus choose nodes with the exact same probabilities.

So an alternative way of performing a step of Price's model is to do the following:

With probability $c/(c + a)$ choose a node in strict proportion to in-degree. Otherwise, choose a node uniformly at random from the set of all nodes.

The choice between the two parts can be achieved, for example, by generating a random number r uniformly in the range $0 \leq r < 1$. If $r < c/(c + a)$ then we choose a node in proportion to in-degree. Otherwise we choose a node uniformly.

Choosing a node uniformly is easily accomplished. Choosing a node in proportion to in-degree is only slightly harder. It can be done rapidly by noting that choosing in proportion to in-degree is equivalent to picking an edge in the network uniformly at random and choosing the node which that edge points to. By definition this makes a node with in-degree q exactly q times as likely to be chosen as a node with in-degree 1, since it has q opportunities to be chosen, one for each of the edges that point to it.

To turn this observation into a computer algorithm we make a list, stored for instance in an ordinary array, of the target of each directed edge in the network. That is, the list's elements contain the node labels i of the nodes to which each edge points. Figure 13.1 shows an example for a small network. Note that the edges do not have to be in any particular order. Any order will do. Nor does the size of the array used to store the list have to match the length of the list exactly; it can contain empty elements at the end as shown in the figure. Indeed, since making already existing arrays larger is difficult in most computer languages, it makes sense to initially create an array that is large enough to hold the longest list we will need. (If the out-degree of nodes is constant this means that it should have length nc , where n is the final number of nodes in the network at the end of the simulation. If out-degree is allowed to fluctuate, then there is a chance the list might grow to be a bit longer than nc , in which case one might create an array of size nc plus a few percent, to be on the safe side.)

Once we have our list, choosing a node in proportion to its in-degree becomes a trivial operation: we simply pick an element uniformly at random

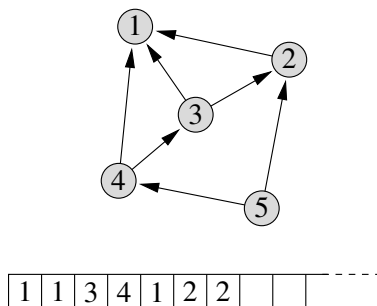


Figure 13.1: The node label list used in the simulation of Price’s model. The list (bottom) contains one entry for the target of each edge in the network (top). In this example, there are three edges that point to node 1 and hence there are three elements containing the number 1 in the list. Similarly there are two containing the number 2, because node 2 is the target of two edges. And so forth.

from the list and choose the node whose label is stored in that element. When a new edge is added to the network, we must also update the list by adding the target of that edge to the end of the list.

Thus our algorithm for creating a new edge is the following:

1. Generate a random number r uniformly in the range $0 \leq r < 1$.
2. If $r < c/(c + a)$, choose an element uniformly at random from the list of targets.
3. Otherwise, choose a node uniformly at random from the set of all nodes.
4. Create an edge linking to the node thus selected, and add that node to the end of the list of targets.

Each step in this process can be accomplished in constant time and hence the creation of all m edges in a network can be accomplished in time $O(m)$. Allowing for setup time and operations required for the creation of each new node, total running time is $O(m + n)$, fast enough to allow the growth of networks with millions of nodes or more.

Figure 13.2a shows the degree distribution of a 100-million-node network generated computationally in this fashion, and the power-law form in the tail of the distribution is clearly visible. A practical problem, however, is the noise in the tail of the histogram, which makes the exact form of the distribution hard to gauge. This is the same problem we encountered for real-world data in Section 10.4.1: the bins in the tail of the histogram have relatively few samples in them, so the statistical fluctuations are large as a fraction of the number of samples. Indeed, simulation data often behave similarly to experimental data

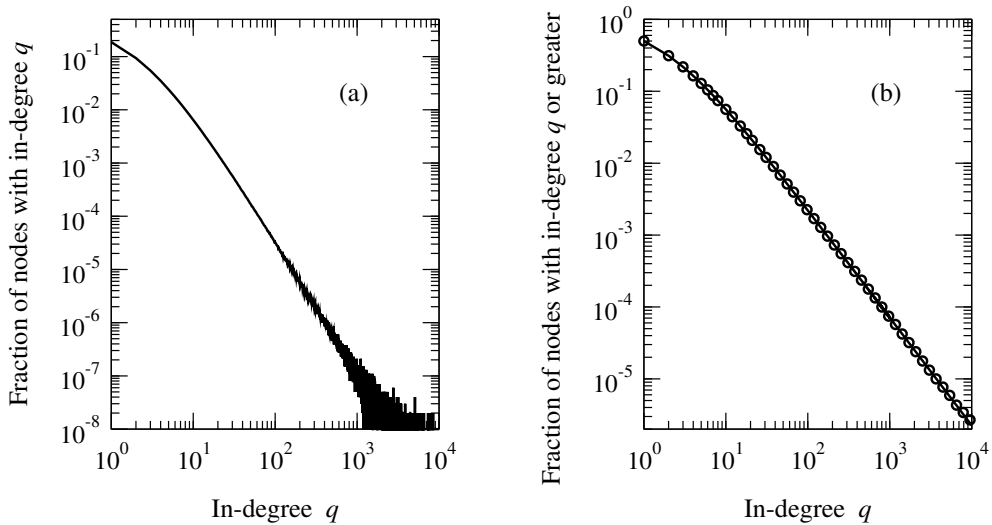


Figure 13.2: Degree distribution in Price's model of a growing network. (a) A histogram of the in-degree distribution for a computer-generated network with $c = 3$ and $a = 1.5$ which was grown until it had $n = 10^8$ nodes. Creation of the network took about 80 seconds on the author's computer using the fast algorithm described in the text. (b) The cumulative distribution function for the same network. The points are the results from the simulation and the solid line is the analytic solution, Eq. (13.34).

in many respects and can often be treated using the same techniques. In this case we can take a hint from Section 10.4.1 and plot a cumulative distribution function instead of a histogram. Recall that the cumulative distribution function P_q is

$$P_q = \sum_{q'=q}^{\infty} p_{q'} \quad (13.32)$$

(see Eq. (10.7)) and is expected to have a power-law tail with an exponent $\alpha - 1 = 1 + a/c$, one less than the exponent of the degree distribution itself. Figure 13.2b shows the cumulative distribution of degrees for our simulation and we now see a much cleaner power-law behavior over several decades in q .

In fact, we can calculate exactly the form the cumulative distribution function should take for Price's model using our analytic solution of the model. To

do this, we make use of the standard integral form for the beta function:⁷

$$B(x, y) = \int_0^1 u^{x-1}(1-u)^{y-1} du. \quad (13.33)$$

Using this expression we find that

$$\begin{aligned} P_q &= \sum_{q'=q}^{\infty} p_{q'} = \frac{1}{B(a, 1+a/c)} \sum_{q'=q}^{\infty} \int_0^1 u^{q'+a-1}(1-u)^{1+a/c} du \\ &= \frac{1}{B(a, 1+a/c)} \int_0^1 u^{a-1}(1-u)^{1+a/c} \sum_{q'=q}^{\infty} u^{q'} du \\ &= \frac{1}{B(a, 1+a/c)} \int_0^1 u^{q+a-1}(1-u)^{a/c} du \\ &= \frac{B(q+a, 1+a/c)}{B(a, 1+a/c)}. \end{aligned} \quad (13.34)$$

Given that $B(x, y)$ goes as x^{-y} for large x (Eq. (13.25)), this implies that indeed the cumulative distribution function has a power-law tail with exponent $1+a/c$.

In Fig. 13.2b we show Eq. (13.34) along with the simulation data, and the simulation and analytic solution agree well, as we would hope.

13.2 THE MODEL OF BARABÁSI AND ALBERT

Price's model of a growing network is an elegant one and the existence of an exact solution showing that its degree distribution has a power-law tail makes a persuasive case for preferential attachment as a possible origin for power-law behavior. Until recently, however, Price's work in this area was not well known outside of the information science community. Preferential attachment

⁷The integral form can be derived by making use of the definition of $B(x, y)$ in terms of gamma functions and the integral form of the gamma function, Eq. (13.15):

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)} = \frac{1}{\Gamma(x+y)} \int_0^{\infty} s^{x-1} e^{-s} ds \int_0^{\infty} t^{y-1} e^{-t} dt.$$

We change variables to $u = s/(s+t)$, $v = s+t$, which gives $s = uv$, $t = (1-u)v$, and a Jacobian of v . Then

$$\begin{aligned} B(x, y) &= \frac{1}{\Gamma(x+y)} \int_0^1 du \int_0^{\infty} v dv (uv)^{x-1} e^{-uv} [(1-u)v]^{y-1} e^{-(1-u)v} \\ &= \frac{1}{\Gamma(x+y)} \int_0^{\infty} v^{x+y-1} e^{-v} dv \int_0^1 u^{x-1} (1-u)^{y-1} du. \end{aligned}$$

The first integral, however, is equal to $\Gamma(x+y)$ by Eq. (13.15) and hence we recover Eq. (13.33).

did not become widely accepted as a mechanism for generating power laws in networks until much later, in the 1990s, when it was independently discovered by Barabási and Albert [40], who proposed their own model of a growing network (along with the name “preferential attachment”). The Barabási–Albert model, which is certainly the best known preferential attachment model in use today, is similar to Price’s, though not identical, being a model of an undirected rather than a directed network.

In the model of Barabási and Albert, nodes are again added one by one to a growing network and each node connects to a suitably chosen set of previously existing nodes. The connections, however, are now undirected and for the model to be solvable the number of connections made by each node must be exactly c (unlike Price’s model, where the number of connections was required only to take an average value of c but might vary from step to step). Note that this implies that c must be an integer, since a node cannot have non-integer degree. Connections are made to nodes with probability precisely proportional to the nodes’ current degree. There is no in- or out-degree now because the network is undirected; connections are made simply in proportion to the (undirected) degree. We will denote the degree of node i by k_i to distinguish it from the directed in-degree q_i of the previous section. As before, nodes and edges are only ever added to the network and never taken away, which means, among other things, that there are no nodes with degree $k < c$. The smallest degree in the network is $k = c$.

One can write down a solution for the model of Barabási and Albert using a master equation method similar to that of Section 13.1 [148,280], but in fact there is no need, because it is straightforward to show that the model is equivalent to a special case of Price’s model. Imagine that, purely for the purposes of our discussion, we give each edge we add to the network a direction, running from the node with which the edge was added to the previously existing node that it connects to. That is, each edge runs from the more recent of the two nodes it joins to the less recent. In this way we convert our network into a directed network in which each node has out-degree exactly c (since this is the number of outgoing edges a node starts with and it never gains any more). And the total degree k_i of a node in the sense of the original undirected network is the sum of the node’s in-degree and out-degree, which is $k_i = q_i + c$ where q_i is the in-degree as before.

But given that the probability of an edge attaching to a node is simply proportional to k_i , it is then also proportional to $q_i + c$, which is the same as Price’s model if we make the particular choice $a = c$. Thus the distribution of in-degrees in this directed network is the same as for Price’s model with $a = c$,

which we find from Eq. (13.21) to be

$$p_q = \frac{B(q+c, 3)}{B(c, 2)}, \quad (13.35)$$

where $B(x, y)$ is once again the Euler beta function, Eq. (13.19). To get the distribution of the total degree we then simply replace $q+c$ by k to get

$$p_k = \begin{cases} \frac{B(k, 3)}{B(c, 2)} & \text{for } k \geq c, \\ 0 & \text{for } k < c. \end{cases} \quad (13.36)$$

This expression can be simplified further by making use of the definition of the beta function in Eq. (13.19) to write

$$p_k = \frac{\Gamma(k)}{\Gamma(k+3)} \frac{\Gamma(c+2)}{\Gamma(c)} \frac{\Gamma(3)}{\Gamma(2)} = \frac{2c(c+1)}{k(k+1)(k+2)} \quad (13.37)$$

for $k \geq c$, where we have used (13.17) in the second equality. In the limit where k becomes large, this gives

$$p_k \sim k^{-3}, \quad (13.38)$$

and hence the Barabási–Albert model generates a degree distribution with a power-law tail and exponent $\alpha = 3$.

Equation (13.37) was first derived by Krapivsky *et al.* [280] and independently by Dorogovtsev *et al.* [148]. A more detailed treatment was later given by Bollobás *et al.* [72], which clarifies precisely the domain of validity of the solution and the possible deviations from the expected value of p_k .

The model of Barabási and Albert can be simulated efficiently on a computer by exploiting the same mapping to Price’s model and the simulation method described in Section 13.1.2. Again we regard the network as a directed one and maintain a list of targets of every directed edge, i.e., the node that each edge points to. Then, setting $a = c$, the algorithm of Section 13.1.2 becomes particularly simple: with probability $\frac{1}{2}$ we choose an element from our list uniformly at random and take the contents of that element as our target node. Otherwise, we choose a target uniformly at random from the set of all nodes currently in existence. Then we create a new edge from the node just added to the selected target and also add that target to the end of our list.

The Barabási–Albert model is attractive for its simplicity—it doesn’t require the offset parameter a of Price’s model and hence has one less parameter to worry about. It is also satisfying that one can write the degree distribution, Eq. (13.37), without using special functions such as the beta and gamma functions that appear in the solution of Price’s model. The price one pays for this simplicity is that the model can no longer match the exponents observed in real networks, being restricted to just a single exponent value $\alpha = 3$.

13.3 TIME EVOLUTION OF THE NETWORK AND THE FIRST MOVER EFFECT

So far we have looked at only one property of the networks generated by preferential attachment models, their overall degree distribution. However, the networks have a number of other interesting features that become apparent when we look a little closer. In particular, since older nodes in the network—those added earlier in the growth process—have more time to acquire links from other nodes, we might expect that they would on average have higher degree. This indeed turns out to be the case, and moreover the effect is a large one. One might imagine that a node twice as old would have twice as many edges on average, but because of the amplification produced by the preferential attachment process the effect turns out to be much larger than this and in fact the oldest nodes in the network end up receiving the lion’s share of all connections. This type of behavior is known as a *first mover effect* or *first mover advantage*, and we can quantify it by calculating the degree distribution of the network as a function of the time at which nodes are created.

Consider a network grown according to Price’s model of Section 13.1. Let $p_q(t, n)$ be the average fraction of nodes in our directed network that were created at time t and have in-degree q when the network has n nodes in total. The time of creation is measured in terms of the number of nodes, the first node having $t = 1$ and the last having $t = n$. Alternatively, you can just think of t as counting the nodes from 1 to n , recording the order in which they were added. Strictly, t need not reflect actual time, because the nodes need not have been added at a constant rate, but if we know the real times at which nodes were added we can easily convert between our timescale and real time.

We can write down a master equation for the evolution of $p_q(t, n)$ as follows. Upon the addition of a new node to the network, the expected number of new edges acquired by previously existing nodes with in-degree q is independent of the time of those nodes’ creation. So, following Eq. (13.2), the expected number acquired by nodes with in-degree q created at time t is

$$np_q(t, n) \times c \times \frac{q + a}{n(c + a)} = \frac{c(q + a)}{c + a} p_q(t, n), \quad (13.39)$$

with the parameters c and a defined as in Section 13.1. Then the master equation takes the form

$$(n+1)p_q(t, n+1) = np_q(t, n) + \frac{c}{c + a} [(q-1+a)p_{q-1}(t, n) - (q+a)p_q(t, n)]. \quad (13.40)$$

If we adopt the convention that $p_q(t, n) = 0$ when $t > n$ for all q , then this equation also gives us the correct result $p_q(n + 1, n + 1) = 0$ when $t = n + 1$.

The only exception to Eq. (13.40) is, as before, for the case $q = 0$, where we get

$$(n + 1)p_0(t, n + 1) = np_0(t, n) + \delta_{t, n+1} - \frac{ca}{c + a}p_0(t, n). \quad (13.41)$$

Note the Kronecker delta, which adds a single node of in-degree zero if $t = n + 1$, giving the correct result $(n + 1)p_0(n + 1, n + 1) = 1$.

Equations (13.40) and (13.41), though correct, don't make much sense in the limit of large n , since the fraction of nodes created at time t goes to zero in this limit because only one node is created at any particular t . So instead we change variables to a rescaled time

$$\tau = \frac{t}{n}, \quad (13.42)$$

which takes values between zero (the oldest nodes) and one (the youngest nodes). At the same time we also change from $p_q(t, n)$ to a density function $\pi_q(\tau, n)$ such that $\pi_q(\tau, n) d\tau$ is the fraction of nodes that have in-degree q and fall in the interval from τ to $\tau + d\tau$. The number of nodes in the interval $d\tau$ is $n d\tau$, which implies that $\pi_q d\tau = p_q \times n d\tau$ and hence

$$\pi_q(\tau, n) = np_q(t, n). \quad (13.43)$$

Being a density function, π_q does not vanish as $n \rightarrow \infty$.

The downside of this variable change is that τ is no longer constant for a given node. A node created at time t has rescaled time t/n when there are n nodes in the network but $t/(n + 1)$ when there are $n + 1$. Thus, in terms of τ and π_q , Eq. (13.40) becomes

$$\begin{aligned} \pi_q\left(\frac{n}{n+1}\tau, n+1\right) &= \pi_q(\tau, n) \\ &+ \frac{c}{c+a} \left[(q-1+a) \frac{\pi_{q-1}(\tau, n)}{n} - (q+a) \frac{\pi_q(\tau, n)}{n} \right]. \end{aligned} \quad (13.44)$$

Now we consider the limit where $n \rightarrow \infty$. If we define the shorthand notation $\pi_q(\tau) = \pi_q(\tau, \infty)$ and the small quantity $\epsilon = 1/n$, Eq. (13.44) becomes

$$\frac{\pi_q(\tau) - \pi_q(\tau - \epsilon\tau)}{\epsilon} + \frac{c}{c+a} [(q-1+a)\pi_{q-1}(\tau) - (q+a)\pi_q(\tau)] = 0, \quad (13.45)$$

where we have dropped terms of order ϵ^2 .

As $n \rightarrow \infty$, we have $\epsilon \rightarrow 0$ and the first two terms become a derivative thus:

$$\lim_{\epsilon \rightarrow 0} \frac{\pi_q(\tau) - \pi_q(\tau - \epsilon\tau)}{\epsilon} = \tau \frac{d\pi_q}{d\tau}, \quad (13.46)$$

and so our master equation becomes a differential equation in this case:

$$\tau \frac{d\pi_q}{d\tau} + \frac{c}{c+a} [(q-1+a)\pi_{q-1}(\tau) - (q+a)\pi_q(\tau)] = 0, \quad (13.47)$$

for $\tau < 1$. For the special case $\tau = 1$, Eq. (13.40) gives $(n+1)p_q(n+1, n+1) = 0$, or $\pi_q(1) = 0$ in the language of our rescaled variables, which gives us a boundary condition on $\pi_q(\tau)$.

Applying similar arguments for the special case $q = 0$, Eq. (13.41) also becomes the differential equation

$$\tau \frac{d\pi_0}{d\tau} - \frac{ca}{c+a} \pi_0(\tau) = 0, \quad (13.48)$$

for $\tau < 1$. For $\tau = 1$, Eq. (13.41) says that $(n+1)p_0(n+1, n+1) = 1$ or equivalently $\pi_0(1) = 1$, which again gives us a boundary condition.⁸

We can solve Eqs. (13.47) and (13.48) by starting with a solution for $q = 0$ and working up through increasing values of q . This is similar to our solution for the degree distribution, Eq. (13.21), except that the equations we are solving are now differential equations. The solution for the $q = 0$ case is straightforward—Eq. (13.48) is homogeneous in π_0 and can be solved by standard methods. You can easily verify that the solution is $\pi_0(\tau) = A\tau^{ca/(c+a)}$ where A is an integration constant. The constant is fixed by the boundary condition $\pi_0(1) = 1$, which implies that $A = 1$ and hence

$$\pi_0(\tau) = \tau^{ca/(c+a)}. \quad (13.49)$$

As a check, we can integrate over τ to get the total fraction of nodes with in-degree zero:

$$\int_0^1 \tau^{ca/(c+a)} d\tau = \frac{1+a/c}{a+1+a/c}, \quad (13.50)$$

which agrees with our previous result for the same quantity, Eq. (13.9).

Now we can use this solution to find $\pi_1(\tau)$. Equation (13.47) tells us that

$$\tau \frac{d\pi_1}{d\tau} - \frac{c}{c+a} (a+1)\pi_1(\tau) = -\frac{c}{c+a} a\pi_0(\tau) = -\frac{ca}{c+a} \tau^{ca/(c+a)}. \quad (13.51)$$

This is again just an ordinary first-order differential equation, although an inhomogeneous one this time (i.e., it has a driving term on the right-hand

⁸Physically, this result arises because there is one node in the time interval between $\tau = 1$ and $\tau = 1 - 1/n$ and it always has in-degree zero. One node is a fraction $1/n$ of the whole network, so we have a density $\pi_0(1) = 1$.

side). We can tackle it in standard fashion. First, we find the general solution for the homogeneous equation in which the right-hand side is set to zero, which is $B\tau^{c(a+1)/(c+a)}$ where B is an integration constant. Then we find any (non-general) solution to the full equation with the driving term included—the obvious one is $a\tau^{ca/(c+a)}$ —and sum the two. The constant is fixed by the boundary condition $\pi_1(1) = 0$, which implies that $B = -a$, and we get

$$\pi_1(\tau) = a\tau^{ca/(c+a)}(1 - \tau^{c/(c+a)}). \quad (13.52)$$

Now, by a similar method, we can use *this* solution to solve for $\pi_2(\tau)$, and so forth to higher and higher values of q . The algebra is tedious, but with persistence you can show that the next two results are

$$\pi_2(\tau) = \frac{1}{2}a(a+1)\tau^{ca/(c+a)}(1 - \tau^{c/(c+a)})^2, \quad (13.53)$$

$$\pi_3(\tau) = \frac{1}{6}a(a+1)(a+2)\tau^{ca/(c+a)}(1 - \tau^{c/(c+a)})^3. \quad (13.54)$$

These results suggest the general solution, first given by Dorogovtsev *et al.* [148]:

$$\begin{aligned} \pi_q(\tau) &= \frac{1}{q!} [a(a+1) \dots (a+q-1)] \tau^{ca/(c+a)} (1 - \tau^{c/(c+a)})^q \\ &= \frac{\Gamma(q+a)}{\Gamma(q+1)\Gamma(a)} \tau^{ca/(c+a)} (1 - \tau^{c/(c+a)})^q, \end{aligned} \quad (13.55)$$

where we have made use again of the convenient property of the gamma function derived in Eq. (13.17) as well as the result that $\Gamma(n+1) = n!$ when n is a positive integer.⁹ With a little work you can verify that this is indeed a complete solution of Eq. (13.47) for all q . As a check we can also integrate over τ to find the total fraction of nodes with in-degree q and confirm that the result agrees with Eq. (13.21). We leave these calculations as an exercise for the reader.¹⁰

Let us take a moment to examine the structure of our solution for $\pi_q(\tau)$ and see what it tells us about the network. The general shape of the solution is shown in Fig. 13.3. Panel (a) shows the distribution of creation times τ for nodes of given in-degree q for various values of q , and for each value there is a clear peak in the distribution, indicating that nodes of a given degree are concentrated around a particular era in the growth of the network. As degree increases, that era gets earlier, so that the times of creation of nodes that ultimately achieve high in-degree are strongly concentrated around the beginning of the growth process.

⁹To prove this we set $x = 1$ in Eq. (13.17) to get $\Gamma(n+1)/\Gamma(1) = n(n-1) \dots 1 = n!$ and from Eq. (13.15) we have $\Gamma(1) = \int_0^\infty e^{-t} dt = 1$.

¹⁰See Exercise 13.10 on page 491.

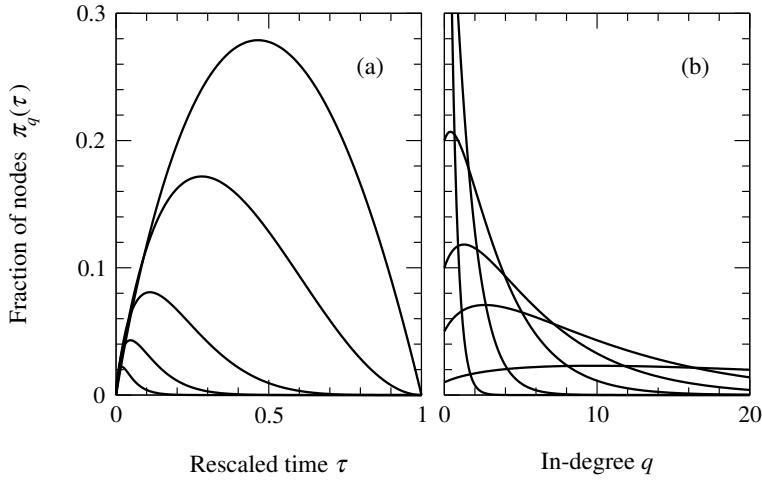


Figure 13.3: Distribution of nodes in Price's model as a function of in-degree and time of creation. The distribution $\pi_q(\tau)$, Eq. (13.55), for $c = 3$ and $a = 1.5$ as (a) a function of τ for (top to bottom) $q = 1, 2, 5, 10$, and 20 , and (b) a function of q for $\tau = 0.01$ (flattest curve), $0.05, 0.1, 0.5$, and 0.9 (steepest curve).

Panel (b) of Fig. 13.3 shows the distribution of in-degrees for nodes created at a selection of different times τ . This distribution also has a peak, then falls off sharply as q becomes large.¹¹ Indeed the distribution falls off roughly exponentially as q becomes large, as we can see from Eq. (13.55) by writing

$$\frac{\Gamma(q+a)}{\Gamma(q+1)\Gamma(a)} = \frac{\Gamma(q+a)}{q\Gamma(q)\Gamma(a)} = \frac{1}{qB(q,a)}, \quad (13.56)$$

where we have again used Euler's beta function, Eq. (13.19). As shown in Section 13.1, the beta function has a power-law tail $B(x, y) \sim x^{-y}$ for large x (Eq. (13.25)) so π_q varies with q as

$$\pi_q(\tau) \sim q^{a-1} (1 - \tau^{c/(c+a)})^q. \quad (13.57)$$

In other words it decays exponentially except for a leading algebraic factor. Thus the degree distribution for nodes with specific values of τ does not follow a power law. The power-law behavior seen in the full degree distribution of the

¹¹Actually, the peak only exists for small values of τ and disappears once τ becomes large enough. There are no peaks in the degree distribution for the $\tau = 0.5$ and $\tau = 0.9$ curves in Fig. 13.3b.

model, Eq. (13.21), only appears when we integrate over all times τ . However, the decay of the exponential in Eq. (13.57) is slower for smaller τ , so older nodes are more likely to have high in-degree than younger ones, as we see in Fig. 13.3.

To investigate this last point further, we can calculate the mean in-degree $\gamma(\tau)$ for a node created at time τ thus:

$$\gamma(\tau) = \sum_{q=0}^{\infty} q \pi_q(\tau) = a(\tau^{-c/(c+a)} - 1). \quad (13.58)$$

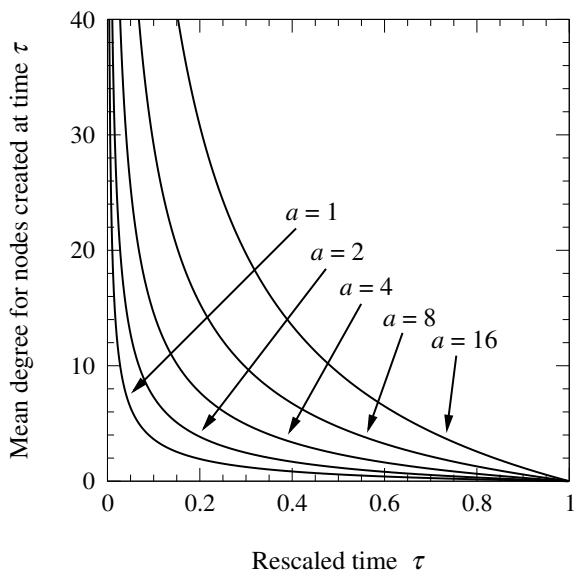


Figure 13.4: Average in-degree of nodes as a function of their time of creation. The average in-degree of nodes in Price’s network model as a function of the rescaled time $\tau = t/n$ at which they were added to the network, in the limit of large n for various values of the parameter a . The out-degree parameter c was in each case $c = 2a$, so that the exponent of the power-law degree distribution $\alpha = 2 + a/c$ (Eq. (13.27)) is 2.5 for all curves, which is a typical value for real-world networks.

Figure 13.4 shows the shape of $\gamma(\tau)$ for a variety of choices of parameters and, as we can see, the mean value of the in-degree always increases with decreasing τ and eventually diverges as τ approaches zero. Note, though, that no node ever actually has $\tau = 0$. The first node is added to the network at time $t = 1$, so the smallest value of τ is $1/n$. Nonetheless, we see that nodes added to the network early have an enormous advantage in terms of in-degree over those added even a little later. For a citation network, for instance, this suggests that the early papers in a field will receive substantially more citations than later ones, purely because they were published first.

This is the first mover effect. First mover effects are seen in many different areas, not just in networks: in any situation where success begets more success, first movers are expected to have a large advantage over others. Any small lead gained early in the process is quickly amplified by the preferential attachment process into a bigger lead and soon the lucky first movers find themselves racing ahead of the pack. Those who enter the game later may also by good luck find themselves with a small lead over their peers, but since there are probably many others already well ahead of them, that lead is not amplified significantly because most of the wealth is already going to the first movers under the preferential attachment rule.

A nice demonstration of this process, although not in the field of networks, was given by Salganik *et al.* [420], who examined the behavior of a group of people downloading popular music online. Salganik *et al.* created a website on

which participants could download and listen to songs by little-known artists for free. Participants were told how many times each song had previously been downloaded and Salganik and co-workers found that there was a clear preferential attachment effect: songs with many previous downloads were downloaded far more than those with few. As a result there was a strong first mover advantage, with songs that took an early lead benefiting from the preferential attachment and turning that lead into a much larger one, resulting in a roughly power-law distribution of the numbers of downloads.

To test the theory that they were seeing a preferential attachment process rather than actual differences in song quality leading to different download rates, Salganik *et al.* then changed the download numbers reported for each song, deliberately misrepresenting the number of times each had been downloaded. They discovered when they did this that the songs with the highest *reported* numbers of downloads were still downloaded most often, even though the reported numbers no longer corresponded to true popularity.¹² These results strongly suggest that success is, at least in this context and at least in part, a result of previous success and that a good way to be successful is to get in at the beginning and get an early lead. Of course, that may be easier said than done. Many people would like to get in at the beginning of a new field of scientific research or a new business opportunity, but it's not always clear how one should do it.

Returning to our network growth model, it is also interesting to ask how the expected in-degree of a node varies with its age after it enters the network. This differs from the expected degree for a particular τ given by Eq. (13.58) because nodes do not have a fixed value of τ . The value of $\tau = t/n$ for a node decreases as time passes because n is increasing. For this reason the behavior of individual nodes is more easily understood in terms of our original non-rescaled time t , which remains constant.

So let t again be the time at which a node is added to the network and let s be the subsequent elapsed time, i.e., the age of the node. Necessarily we have $t + s = n$ and hence

$$\tau = \frac{t}{n} = \frac{t}{t + s}. \quad (13.59)$$

Substituting this expression into Eq. (13.58), we then find the expected in-

¹²Salganik *et al.* did find a weak effect of song quality—songs that had proved popular when the download numbers were reported faithfully continued to do better than expected even when the download numbers were changed.

degree $\gamma_t(s)$ of the node added at time t , as a function of its age s , to be

$$\gamma_t(s) = a \left[\left(1 + \frac{s}{t} \right)^{c/(c+a)} - 1 \right]. \quad (13.60)$$

When a node is first added to the network and $s \ll t$, we can expand in the small quantity s/t to get

$$\gamma_t(s) \approx \frac{ca}{c+a} \left(\frac{s}{t} \right). \quad (13.61)$$

In other words, the in-degree of a node initially grows linearly with the age of the node, on average, but with a constant of proportionality that is smaller the later the node entered the network—again we see that there is a substantial advantage for nodes that enter early.

As the node ages, there is a crossover to another regime around the point $s = t$, i.e., at the point where the node switches from being in the younger half of the population to being in the older. For $s \gg t$, we have

$$\gamma_t(s) \approx a \left(\frac{s}{t} \right)^{c/(c+a)}, \quad (13.62)$$

so the growth is slower than linear for older nodes but still favors nodes that appear early. Figure 13.5 shows the behavior of $\gamma_t(s)$ with time for nodes created at a selection of different times t .

All of these results can be applied to the Barabási–Albert model as well by setting $a = c$ with c an integer and writing the formulas in terms of total degree $k = q + c$ rather than in-degree. For instance, the joint degree/time distribution, Eq. (13.55), becomes

$$\pi_k(\tau) = \binom{k-1}{c-1} (\sqrt{\tau})^c (1 - \sqrt{\tau})^{k-c} \quad (13.63)$$

for $k \geq c$ and $\pi_k(\tau) = 0$ for $k < c$. This result was first given by Krapivsky and Redner [279] for the case $c = 1$.

13.4 EXTENSIONS OF PREFERENTIAL ATTACHMENT MODELS

Many extensions and generalizations of preferential attachment models have been suggested, typically addressing questions about what happens when we vary the details of the model definition or attempt to make the model more faithful to the way real networks grow. For example, by contrast with citations, links on the Web can be added not just at the moment a node is created but at any later time too. And links between web pages can disappear at any time, as indeed can the pages themselves. There is also no obvious reason why preferential attachment processes need be linear in the degree. What happens if

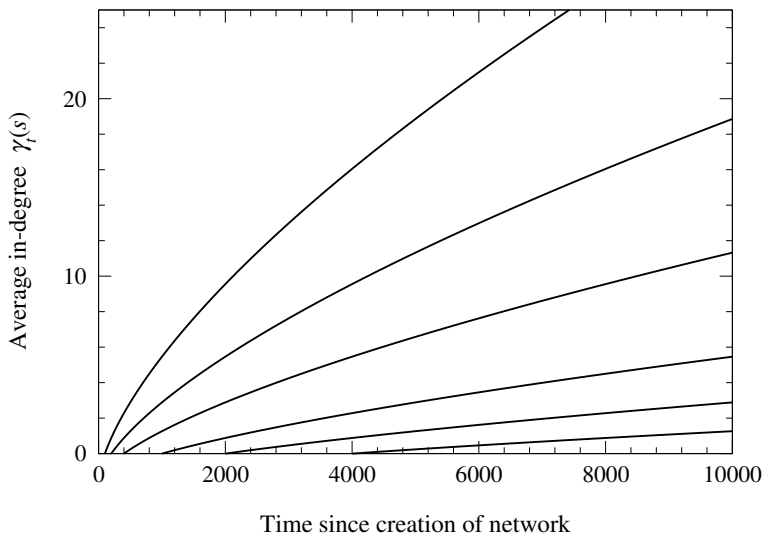


Figure 13.5: Average in-degree of nodes created at different times. The curves show the average in-degrees in Price’s model of nodes created at times (top to bottom) $t = 100, 200, 400, 1000, 2000,$ and 4000 as a function of time since the creation of the network. The model parameters were $c = 3$ and $a = 1.5$.

they are non-linear? In this section we describe modified preferential attachment models that address each of these issues. In the interest of simplicity, we describe the developments in the context of the Barabási–Albert model, rather than the more general Price model. Generalizations of Price’s model are certainly possible but the algebra is in many cases unwieldy, and the main conclusions are easier to understand in the context of the simpler model.

13.4.1 ADDITION OF EXTRA EDGES

Price proposed his model of a growing network with citation networks in mind. Since the bibliography of a paper cannot be changed after the paper is published, the edges in a citation network are effectively frozen in place from the moment they are first created, and Price’s model mimics this behavior with edges being added only at the moment a node is created and never moved or removed thereafter.

This is not true of all networks, however. The World Wide Web, for example, is constantly changing. Links between web pages can, and often are, added or removed after the pages are created. This state of flux is not captured by

Price's model or by the Barabási–Albert model of Section 13.2. Yet the Web still has a power-law degree distribution. This leads us to wonder whether we can make a generalized model that includes the addition and removal of edges after nodes are created but still generates power-law distributions. It turns out that we can, as we now describe.

We first consider the relatively simple case, studied by a number of authors [14, 145, 281], in which edges can be added between nodes after the nodes are created, but no edges are ever taken away. The case of edge removal is more complex and is considered in the following section. The model we consider here is a generalization of the Barabási–Albert model in which nodes are added to the network one by one as before and each starts out with c undirected edges which attach to other nodes with probability proportional to degree k . But we now include a second process in the model as well: at each step some number w of extra edges are added to the network with *both* ends attaching to nodes chosen in proportion to degree. Thus when the network has n nodes it will have a total of $n(c + w)$ edges. (In fact, it is only necessary that an *average* number w of extra edges be added at each step. The actual number can fluctuate around this figure, and the net result, in the limit of large network size, will be the same. This allows us to give w a non-integer value if we wish.)

This model turns out to be quite easy to solve given the results of previous sections. The only difference between it and the standard Barabási–Albert model is that, instead of c new ends of edges attaching to old nodes for every new node added, we now have $c + 2w$ new ends of edges—two extra for each of the w extra edges. The probability of attachment of any one of those ends of edges to a particular node i is $k_i / \sum_i k_i$. The sum in the denominator of this expression is equal to twice the number of edges in the network (see Eq. (6.13)), so $\sum_i k_i = 2n(c + w)$.

Then, if $p_k(n)$ denotes the fraction of nodes with degree k when the network has n nodes in total, the number of nodes of degree k receiving a new edge, when one node is added to the network, is

$$np_k(n) \times (c + 2w) \times \frac{k}{2n(c + w)} = \frac{c + 2w}{2(c + w)} kp_k(n). \quad (13.64)$$

We can use this result to write a master equation for $p_k(n)$ thus:

$$(n + 1)p_k(n + 1) = np_k(n) + \frac{c + 2w}{2(c + w)} [(k - 1)p_{k-1}(n) - kp_k(n)], \quad (13.65)$$

for $k > c$ and

$$(n + 1)p_c(n + 1) = np_c(n) + 1 - \frac{c + 2w}{2(c + w)} cp_c(n), \quad (13.66)$$

for $k = c$. (There are, as before, no nodes of degree less than c .) Taking the limit of large n and writing $p_k = p_k(\infty)$, these equations simplify to

$$p_k = \frac{c + 2w}{2(c + w)} [(k - 1)p_{k-1} - kp_k] \quad \text{for } k > c, \quad (13.67)$$

$$p_c = 1 - \frac{c + 2w}{2(c + w)} cp_c \quad \text{for } k = c. \quad (13.68)$$

Rearranging these expressions along the lines of Eqs. (13.9) to (13.21), we then find that

$$p_k = \frac{B(k, \alpha)}{B(c, \alpha - 1)}, \quad (13.69)$$

where $B(x, y)$ is again the Euler beta function, Eq. (13.19), and

$$\alpha = 2 + \frac{c}{c + 2w}. \quad (13.70)$$

Since $B(x, y)$ goes as x^{-y} for large x (Eq. (13.25)), our degree distribution has a power-law tail with exponent α . For the special case $w = 0$, in which no additional edges are added to the network, we recover the standard result $\alpha = 3$ for the Barabási–Albert model; for $w > 0$ we get exponents in the range $2 < \alpha < 3$, which agrees nicely with the values typically observed for degree distributions on the Web (see Table 10.1). Bear in mind though that the Web is a directed network while the model described here is undirected. If we want to build a model of a directed network we would need to start with something like the Price model of Section 13.1. Generalizations of Price’s model that include the addition of extra edges as in this section are certainly possible—see, for example, Krapivsky *et al.* [281].

13.4.2 REMOVAL OF EDGES

Now consider the case of a network in which edges can be removed. To keep things simple let us first consider the case where edges can be removed at any time but are only added at the initial creation of a node, as in the standard Barabási–Albert model. (We will consider the general case of addition and removal at any time shortly.)

There are many ways in which edges could be removed from a network, but let us consider the most basic case in which they are simply deleted uniformly at random. What then is the probability that a particular node i loses an edge when a single edge is removed from the network? When an edge is deleted both of its two ends vanish. Given that the deletion is uniformly random over edges, the probability that one of those two ends is attached to node i is simply

proportional to the total number of ends attached to i , which is equal to the degree k_i . Properly normalized, the probability that node i loses an edge is thus $2k_i / \sum_i k_i$, the factor of two coming from the two ends of the edge. In other words, the random deletion of edges is like a type of preferential attachment in reverse: the higher the degree of the node, the more likely it is to lose an edge.

So consider the undirected network model in which nodes with degree c are added to the network following the normal preferential attachment scheme and an average of v edges are deleted at random for each node added. (As with the model of Section 13.4.1, the actual number of edges deleted can fluctuate about the mean and v can take a non-integer value if we wish.) To ensure that the number of edges in the network grows, rather than shrinking to zero and vanishing, we require that the net number of edges added per node $c - v$ be positive, i.e., that $v < c$. Then when the network has n nodes the number of edges will be $n(c - v)$ and $\sum_i k_i = 2n(c - v)$.

In writing down a master equation for this model there are several processes we need to consider. As before, the number of nodes with degree k increases whenever a node of degree $k - 1$ gains a new edge and decreases when a node of degree k gains a new edge. By an argument analogous to the one leading to Eq. (13.64), the number of nodes of degree k gaining an edge per node added to the network is

$$np_k(n) \times c \times \frac{k}{2n(c - v)} = \frac{c}{2(c - v)} kp_k(n). \quad (13.71)$$

But we also now have a new process in which a node can lose an edge, which means that the number of nodes of degree k also increases when a node of degree $k + 1$ loses an edge and decreases when a node of degree k loses an edge. The number of nodes of degree k losing an edge per node added is given by

$$np_k(n) \times v \times \frac{2k}{2n(c - v)} = \frac{v}{c - v} kp_k(n). \quad (13.72)$$

Note that, by contrast with the original Barabási–Albert model, nodes can now have any degree $k \geq 0$ —nodes can lose any or all of their edges, right down to the last one, so there is no restriction $k \geq c$ on the degree as before.

Putting together the above results, our master equation now takes the form

$$(n + 1)p_k(n + 1) = np_k(n) + \frac{c}{2(c - v)}(k - 1)p_{k-1}(n) + \frac{v}{c - v}(k + 1)p_{k+1}(n) - \frac{c + 2v}{2(c - v)}kp_k(n) \quad (13.73)$$

for $k \neq c$ and

$$(n+1)p_c(n+1) = np_c(n) + 1 + \frac{c}{2(c-v)}(c-1)p_{c-1}(n) + \frac{v}{c-v}(c+1)p_{c+1}(n) - \frac{c+2v}{2(c-v)}cp_c(n) \quad (13.74)$$

for $k = c$. These two equations can conveniently be combined by writing

$$(n+1)p_k(n+1) = np_k(n) + \delta_{kc} + \frac{c}{2(c-v)}(k-1)p_{k-1}(n) + \frac{2v}{2(c-v)}(k+1)p_{k+1}(n) - \frac{c+2v}{2(c-v)}kp_k(n), \quad (13.75)$$

where δ_{kc} is the Kronecker delta.

The only exception to this master equation is for the case $k = 0$, where the term proportional to $k - 1$ vanishes because there are no nodes of degree -1 . A simple way of enforcing this exception is to define $p_{-1}(n) = 0$ for all n , in which case Eq. (13.75) then applies for all $k \geq 0$. We will adopt this convention henceforth.

The model as we have described it so far incorporates the processes of node addition and edge removal, but, given Eq. (13.75), it is only a small extra step to incorporate the edge addition process of Section 13.4.1 as well. If as before we add w extra edges per node added, then $c + w - v$ edges are added net per node, and our master equation becomes

$$(n+1)p_k(n+1) = np_k(n) + \delta_{kc} + \frac{c+2w}{2(c+w-v)}(k-1)p_{k-1}(n) + \frac{v}{c+w-v}(k+1)p_{k+1}(n) - \frac{c+2w+2v}{2(c+w-v)}kp_k(n). \quad (13.76)$$

The equation for edge removal only, Eq. (13.75), can then be considered a special case of this equation with $w = 0$. As before, we require that the net number of edges added per node be positive, or $v < c + w$.

Taking the limit $n \rightarrow \infty$ and writing $p_k = p_k(\infty)$ we now find that the limiting degree distribution satisfies

$$p_k = \delta_{kc} + \frac{c+2w}{2(c+w-v)}(k-1)p_{k-1} + \frac{v}{c+w-v}(k+1)p_{k+1} - \frac{c+2w+2v}{2(c+w-v)}kp_k. \quad (13.77)$$

This equation differs in a crucial way from the master equations we have encountered previously, such as Eq. (13.7), because the right-hand side contains

terms for nodes of three different degrees—namely $k - 1$, k , and $k + 1$ —rather than just two. This makes the equation substantially more difficult to solve. We can no longer rearrange it to derive an expression for p_k in terms of p_{k-1} and then apply that expression repeatedly to itself. A solution is still possible, but it's not simple. Here we give just an outline of the method. The details, for those interested in them, are spelled out by Moore *et al.* [339].¹³

The basic strategy for solving Eq. (13.77) is to use a generating function of the kind we introduced in Section 12.10. We define

$$g(z) = \sum_{k=0}^{\infty} p_k z^k. \quad (13.78)$$

Substituting for p_k from Eq. (13.77) we get

$$g(z) = \sum_{k=0}^{\infty} \delta_{kc} z^k + \frac{c + 2w}{2(c + w - v)} \sum_{k=0}^{\infty} (k - 1) p_{k-1} z^k + \frac{2v}{2(c + w - v)} \sum_{k=0}^{\infty} (k + 1) p_{k+1} z^k - \frac{c + 2w + 2v}{2(c + w - v)} \sum_{k=0}^{\infty} k p_k z^k. \quad (13.79)$$

The first term on the right is simple—it is equal to z^c . The others are a little more complicated. Consider the second term, for example. Note that the first term in the sum, the term for $k = 0$, is necessarily zero because, as we have said, $p_{-1} = 0$. Hence we can write

$$\begin{aligned} \sum_{k=0}^{\infty} (k - 1) p_{k-1} z^k &= \sum_{k=1}^{\infty} (k - 1) p_{k-1} z^k = \sum_{k=0}^{\infty} k p_k z^{k+1} \\ &= z^2 \sum_{k=0}^{\infty} k p_k z^{k-1} = z^2 \frac{d}{dz} \sum_{k=0}^{\infty} p_k z^k \\ &= z^2 \frac{dg}{dz}, \end{aligned} \quad (13.80)$$

where in the first line we have made the substitution $k - 1 \rightarrow k$ and in the second line we have made use of the fact that the $k = 0$ term is again zero (because of the factor of k).

For the third and fourth terms in (13.79) we can similarly write

$$\sum_{k=0}^{\infty} (k + 1) p_{k+1} z^k = \sum_{k=1}^{\infty} k p_k z^{k-1} = \sum_{k=0}^{\infty} k p_k z^{k-1} = \frac{dg}{dz}, \quad (13.81)$$

¹³In fact, Moore *et al.* give a solution for a model in which nodes rather than edges are deleted, but the two can be treated by virtually the same means. The calculation given here is adapted from their work with only minor changes.

and

$$\sum_{k=0}^{\infty} k p_k z^k = z \sum_{k=0}^{\infty} k p_k z^{k-1} = z \frac{dg}{dz}. \quad (13.82)$$

Combining Eqs. (13.79) to (13.82) and rearranging, we then get

$$\frac{(c+2w)z-2v}{2(c+w-v)}(1-z)\frac{dg}{dz} + g(z) = z^c. \quad (13.83)$$

This is a first-order linear differential equation and is solvable by standard—if tedious—methods. To cut a long story short, one can find an integrating factor for the left-hand side and hence express the solution in terms of an integral that, provided $v < \frac{1}{2}c + w$, can be reduced by repeated integration by parts to

$$p_k = Ak^{-\alpha} \int_0^k \frac{(1-x/k)^k}{(1-\gamma x/k)^k} x^{\alpha-2} dx, \quad (13.84)$$

for $k \geq c$, where A is a k -independent normalizing constant and

$$\alpha = 2 + \frac{v-w}{c+2w-2v}, \quad (13.85)$$

$$\gamma = \frac{2v}{c+2w}. \quad (13.86)$$

The remaining integral can be written in terms of hypergeometric functions, but we can find the asymptotic behavior of the degree distribution for large k more directly by noting that as k becomes large

$$\left(1 - \frac{x}{k}\right)^k \rightarrow e^{-x}, \quad \left(1 - \frac{\gamma x}{k}\right)^k \rightarrow e^{-\gamma x}, \quad (13.87)$$

so that

$$p_k \sim k^{-\alpha} \int_0^{\infty} e^{-(1-\gamma)x} x^{\alpha-2} dx = \frac{\Gamma(\alpha-1)}{(1-\gamma)^{\alpha-1}} k^{-\alpha}. \quad (13.88)$$

Thus, we once again find that our degree distribution has a power-law tail, with an exponent given this time by Eq. (13.85). Note that this exponent can take values both greater than and less than 2. What's more, for the case where $v = \frac{1}{2}c + w$, it actually becomes infinite. Moore *et al.* [339] showed that at this point we lose the power-law behavior and the distribution becomes instead a stretched exponential. Up until this point, however, the distribution follows a power law, albeit with a very large exponent as v grows larger. For values of $v > \frac{1}{2}c + w$, the solution becomes nonsensical, with a negative value of α , and one must return to the original differential equation (13.83) to find the solution for this case. We leave the developments, however, as an exercise for the especially avid reader.

Before we leave this topic let us point out that the methods used to solve Eq. (13.77) can also be used to calculate what happens when we remove not edges but nodes from our network. Loss of nodes does occur in some networks, such as the World Wide Web, so it is potentially of interest to ask what effect it has on the degree distribution. In fact, the solution for this case is very similar to the solution for loss of edges, with a power-law distribution and an exponent that depends on the node loss rate, diverging as the rate of loss approaches the rate at which nodes are added. The details can be found in [339].

13.4.3 NON-LINEAR PREFERENTIAL ATTACHMENT

In the models we have considered so far, the probability that a new edge attaches to a node is linear in the degree of the node. Although this is a reasonable first guess about the way things might work, it's certainly possible that attachment processes might not be linear. Indeed, there is some empirical evidence that this is the case. For instance, Jeong *et al.* [251] looked at the growth of several real-world networks, measuring the rate at which nodes acquired new edges. To avoid problems associated with the fact that the rate can depend not only on degree but also on the total size n of the network (see Eq. (13.1)), they restricted their observations to relatively short intervals of time. The measured rates, plotted as a function of node degree, showed that for some networks there was a roughly linear preferential attachment effect, but for others attachment appeared to be non-linear, with nodes gaining new edges at a rate going as some power γ of the degree with γ being less than 1, typically around $\gamma = 0.8$.

We can study the effects of non-linear preferential attachment by building a model along the lines of those in this chapter. Following an approach introduced by Krapivsky *et al.* [280], we define an *attachment kernel*, denoted a_k , which specifies how the probability of a node gaining an edge depends on its degree. For the model of Barabási and Albert, where attachment is simply proportional to degree, the attachment kernel would be $a_k = k$. For the non-linear attachment observed by Jeong *et al.* it would be $a_k = k^\gamma$. Note that the attachment kernel is not itself a probability, merely a functional form. The correctly normalized probability that a newly added edge attaches to a specific node having degree k is $a_k / \sum_i a_{k_i}$.

Now consider a growing undirected network of the type discussed in previous sections and let $p_k(n)$ be the fraction of nodes with degree k when the network has n nodes. As before, c new edges will be added to the network with each new node, but preferential attachment is now non-linear, governed by the attachment kernel a_k , which means that, by analogy with Eq. (13.2), the expected number of nodes of degree k receiving a new edge when a single new

node is added to the network is

$$np_k(n) \times c \times \frac{a_k}{\sum_i a_{k_i}} = \frac{c}{\mu(n)} a_k p_k(n), \quad (13.89)$$

where $\mu(n)$ is the average value of the attachment kernel over all nodes:

$$\mu(n) = \frac{1}{n} \sum_{i=1}^n a_{k_i} = \sum_k a_k p_k(n). \quad (13.90)$$

Now the master equation for $p_k(n)$ is

$$(n+1)p_k(n+1) = np_k(n) + \frac{c}{\mu(n)} [a_{k-1}p_{k-1}(n) - a_k p_k(n)]. \quad (13.91)$$

As before, the term in $p_{k-1}(n)$ represents new nodes of degree k created when nodes of degree $k-1$ receive new edges and the final term in $p_k(n)$ represents nodes of degree k that are lost when they gain new edges to become nodes of degree $k+1$.

The only exception to this equation is for nodes of degree c , for which

$$(n+1)p_c(n+1) = np_c(n) + 1 - \frac{c}{\mu(n)} a_c p_c(n). \quad (13.92)$$

(And there are no nodes of degree less than c , since all nodes are created with degree c initially and edges are never removed.)

Taking the limit as $n \rightarrow \infty$ and writing $p_k = p_k(\infty)$ and $\mu = \mu(\infty)$, these equations become

$$p_k = \frac{c}{\mu} [a_{k-1}p_{k-1} - a_k p_k] \quad (13.93)$$

for $k > c$ and

$$p_c = 1 - \frac{ca_c}{\mu} p_c. \quad (13.94)$$

Note that μ depends via Eq. (13.90) on the degree distribution, which we don't yet know. For now, however, it will be enough that μ is independent of k ; we will derive an expression for its exact value in a moment.

Equations (13.93) and (13.94) can be rearranged to give

$$p_c = \frac{\mu/c}{a_c + \mu/c} \quad (13.95)$$

and

$$p_k = \frac{a_{k-1}}{a_k + \mu/c} p_{k-1}. \quad (13.96)$$

Applying the latter repeatedly we get

$$\begin{aligned}
 p_k &= \frac{a_{k-1}a_{k-2}\dots a_c}{(a_k + \mu/c)\dots(a_{c+1} + \mu/c)} p_c \\
 &= \frac{\mu}{ca_k} \frac{a_k \dots a_c}{(a_k + \mu/c)\dots(a_c + \mu/c)} \\
 &= \frac{\mu}{ca_k} \prod_{r=c}^k \left[1 + \frac{\mu}{ca_r} \right]^{-1}. \tag{13.97}
 \end{aligned}$$

All we need to complete our solution is the value of μ . Taking Eq. (13.90) and letting $n \rightarrow \infty$, we get

$$\mu = \sum_{k=c}^{\infty} a_k p_k = \frac{\mu}{c} \sum_{k=c}^{\infty} \prod_{r=c}^k \left[1 + \frac{\mu}{ca_r} \right]^{-1}. \tag{13.98}$$

Canceling a factor of μ from both sides we arrive at the equation

$$\sum_{k=c}^{\infty} \prod_{r=c}^k \left[1 + \frac{\mu}{ca_r} \right]^{-1} = c. \tag{13.99}$$

In principle we should be able to solve this equation for μ and substitute the result into Eq. (13.97) to get the complete degree distribution. In practice, unfortunately, the equation is not solvable in closed form for most choices of the attachment kernel a_k , although an approximate value for μ can usually be calculated numerically on a computer.

Even without knowing μ , however, we can still find the overall functional form of p_k , which is enough to answer many of the questions we are interested in. As an example, consider a network of the type observed by Jeong *et al.* [251] and discussed at the start of this section, in which attachment goes as k^γ for some positive constant γ , and let us assume (as found by Jeong *et al.*) that $\gamma < 1$. The solution for this particular choice was given by Krapivsky *et al.* [280] and shows a number of interesting features.

Putting $a_k = k^\gamma$ in Eq. (13.97) gives

$$p_k = \frac{\mu}{ck^\gamma} \prod_{r=c}^k \left[1 + \frac{\mu}{cr^\gamma} \right]^{-1}. \tag{13.100}$$

This degree distribution turns out *not* to have a power-law tail, in contrast to the case of linear preferential attachment. We can see this by writing

$$\prod_{r=c}^k \left[1 + \frac{\mu}{cr^\gamma} \right]^{-1} = \exp \left[- \sum_{r=c}^k \ln \left(1 + \frac{\mu}{cr^\gamma} \right) \right] \tag{13.101}$$

and then expanding the logarithm as a Taylor series in μ/cr^γ :

$$\begin{aligned} \sum_{r=c}^k \ln\left(1 + \frac{\mu}{cr^\gamma}\right) &= - \sum_{r=c}^k \sum_{s=1}^{\infty} \frac{(-1)^s}{s} \left(\frac{\mu}{c}\right)^s r^{-s\gamma} \\ &= - \sum_{s=1}^{\infty} \frac{(-1)^s}{s} \left(\frac{\mu}{c}\right)^s \sum_{r=c}^k r^{-s\gamma}. \end{aligned} \quad (13.102)$$

The sum over r cannot be expressed in closed form, but we can approximate it using the *trapezoidal rule*,¹⁴ which says that for any function $f(r)$:

$$\sum_{r=a}^b f(r) = \int_a^b f(r) dr + \frac{1}{2}[f(a) + f(b)] + O(f'(b) - f'(a)). \quad (13.103)$$

(For those not familiar with it, the derivation of the trapezoidal rule is illustrated in Fig. 13.6.¹⁵)

In our case $f(r) = r^{-s\gamma}$ and Eq. (13.103) gives

$$\sum_{r=c}^k r^{-s\gamma} = A_s + \frac{k^{1-s\gamma}}{1-s\gamma} + \frac{1}{2}k^{-s\gamma} + O(k^{-(s\gamma+1)}), \quad (13.104)$$

where A_s is a constant depending on s (and on c) but not on k .

Consider now what happens when k becomes large. Since γ is positive the term in $k^{-s\gamma}$ and all subsequent terms vanish as $k \rightarrow \infty$ and Eq. (13.102) becomes

$$\sum_{r=c}^k \ln\left(1 + \frac{\mu}{cr^\gamma}\right) \simeq A - \sum_{s=1}^{\infty} \frac{(-1)^s}{s} \left(\frac{\mu}{c}\right)^s \frac{k^{1-s\gamma}}{1-s\gamma}, \quad (13.105)$$

where A is a k -independent constant equal to $\sum_{s=1}^{\infty} A_s (-\mu/c)^s / s$.

This expression can be simplified still further by noting that, in the limit $k \rightarrow \infty$, all terms in $k^{1-s\gamma}$ where $1-s\gamma < 0$ also vanish. Thus for any given value of γ we need keep terms in k up to only a certain value of s . The simplest case is when $\frac{1}{2} < \gamma < 1$. In this case only the term for $s = 1$ grows as k increases, all others vanishing, and

$$\sum_{r=c}^k \ln\left(1 + \frac{\mu}{cr^\gamma}\right) \simeq A + \frac{\mu k^{1-\gamma}}{c(1-\gamma)} \quad (13.106)$$

¹⁴Also called the trapezium rule in British English.

¹⁵Explicit expressions are known for the correction terms (the terms in $f'(a)$ and $f'(b)$)—they are given in terms of the Bernoulli numbers by the so-called Euler–Maclaurin formula [2]—but they're not necessary in our application because these terms vanish anyway.

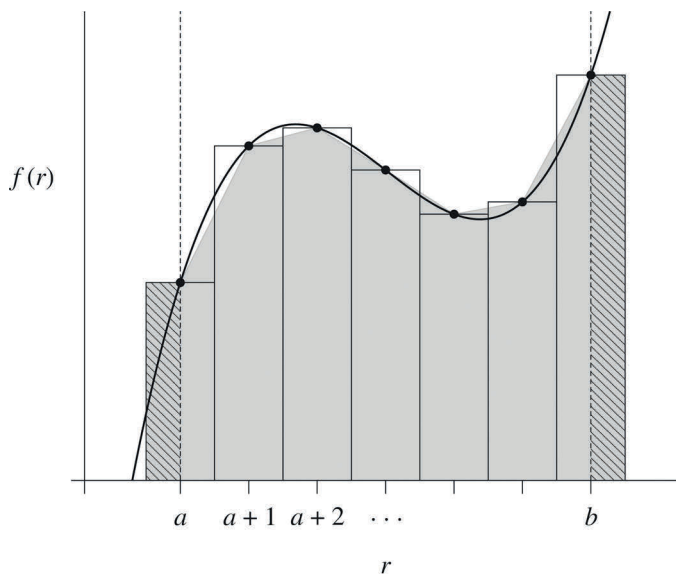


Figure 13.6: The trapezoidal rule. The trapezoidal rule is normally used to approximate an integral by a sum, but here we use it in the reverse direction to approximate a sum by an integral. The sum of the function $f(r)$ from $r = a$ to $r = b$ is equal to the sum of the areas of the rectangular bars, which is also equal to the area shaded in gray. This shaded area can be approximated by the integral of $f(r)$ between a and b (smooth curve) plus the two extra rectangular sections at either end (hatched), which have area $\frac{1}{2}f(a)$ and $\frac{1}{2}f(b)$ respectively. Add everything up and we get Eq. (13.103). The error in the approximation is equal to the sum of the relatively small regions between the curve and the edge of the shaded area.

as $k \rightarrow \infty$.

Now, combining Eqs. (13.100), (13.101), and (13.106), we find that the asymptotic form of p_k is

$$p_k \sim k^{-\gamma} \exp\left(-\frac{\mu k^{1-\gamma}}{c(1-\gamma)}\right), \quad (13.107)$$

for $\frac{1}{2} < \gamma < 1$.

Distributions of this general form, in which the dominant contribution to the probability falls off as the exponential of a power of k , are called *stretched exponentials*. Since the exponent $1 - \gamma$ is less than 1, the distribution falls off more slowly than an ordinary exponential in k , which is why we called it

“stretched.”¹⁶ On the other hand, the distribution still falls off a good deal faster than the power law that we found in the case of linear preferential attachment, and this is really the important point here. This calculation reveals that the power-law distribution in the Barabási–Albert model is a special feature of the linear attachment process assumed by that model. (Note that this observation is valid even though we haven’t calculated the value of the constant μ . The general functional form of the degree distribution doesn’t depend on the value of the constant.)

For other values of γ the calculation is similar but involves more terms in Eq. (13.105). For instance, if $\frac{1}{3} < \gamma < \frac{1}{2}$ then the terms in $k^{1-s\gamma}$ for $s = 1$ and 2 both grow as k becomes large while all others vanish, and we find that

$$\sum_{r=k_0+1}^k \ln\left(1 + \frac{\mu}{cr^\gamma}\right) \simeq A + \frac{\mu k^{1-\gamma}}{c(1-\gamma)} - \frac{\mu^2 k^{1-2\gamma}}{2c^2(1-2\gamma)}, \quad (13.108)$$

which gives

$$p_k \sim k^{-\gamma} \exp\left(-\frac{\mu k^{1-\gamma}}{c(1-\gamma)} + \frac{\mu^2 k^{1-2\gamma}}{2c^2(1-2\gamma)}\right). \quad (13.109)$$

In between the solutions (13.107) and (13.109) there is a special case solution when γ is exactly equal to $\frac{1}{2}$. For $\gamma = \frac{1}{2}$ and $s = 2$ the integral in Eq. (13.103) gives rise not to a power of k but to a log and Eq. (13.102) becomes

$$\sum_{r=c}^k \ln\left(1 + \frac{\mu}{cr^\gamma}\right) \simeq A + \frac{2\mu}{c} \sqrt{k} - \frac{\mu^2}{2c^2} \ln k, \quad (13.110)$$

all other terms vanishing in the limit of large k . Substituting this expression into Eq. (13.101), we then arrive at

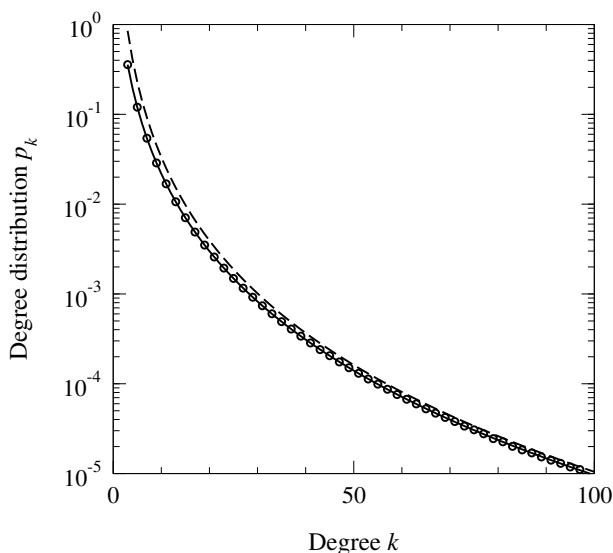
$$p_k \sim (\sqrt{k})^{\mu^2/c^2-1} \exp\left(-\frac{2\mu}{c} \sqrt{k}\right), \quad (13.111)$$

for $\gamma = \frac{1}{2}$.

We can continue in this vein. There are distinct solution forms for $\frac{1}{4} < \gamma < \frac{1}{3}$ and $\frac{1}{5} < \gamma < \frac{1}{4}$ and so forth, as well as special case solutions for $\gamma = \frac{1}{3}, \frac{1}{4}, \frac{1}{5}$, and so forth. Figure 13.7 shows the degree distribution for the case $\gamma = 0.8$, along with the asymptotic form (13.107). Note the convex form of the curve on the semilogarithmic scales, which indicates a function decaying slower than an exponential.

¹⁶Confusingly, however, people often still call it a stretched exponential even when the exponent is greater than 1, although one does occasionally see the latter case referred to (more accurately) as a “compressed exponential.”

Figure 13.7: Degree distribution for sublinear preferential attachment. The fraction p_k of nodes with degree k in a growing network with attachment kernel k^γ as described in the text. In this case $\gamma = 0.8$ and $c = 3$. The points are results from computer simulations, averaged over 100 networks of (final) size 10^7 nodes each. The solid line is the exact solution, Eq. (13.100), evaluated numerically. The dashed line is the asymptotic form, Eq. (13.107), with the overall constant of proportionality chosen to coincide with the exact solution for large k .



One can also calculate the degree distribution for superlinear preferential attachment, i.e., for values of γ greater than 1. This case also shows a non-power-law distribution with some interesting behaviors: it turns out that for $\gamma > 1$ the typical behavior is for one node to emerge as a “leader” in the network, gaining a non-zero fraction of all edges, with the rest of the nodes having small degree (almost all having degree less than some fixed constant). Readers interested in these developments can find them described in detail in Ref. [280].

Many other extensions and variations of preferential attachment models have been studied in addition to those described in this chapter. If you’re interested in learning more, there are a number of review articles that go into the subject in some detail—see Refs. [15, 68, 147]. The rest of this chapter is devoted to the discussion of other models of network formation and growth that don’t rely on preferential attachment.

13.5 NODE COPYING MODELS

Preferential attachment models offer a plausible, if simplified, explanation for power-law degree distributions in networks such as citation networks and the World Wide Web. Preferential attachment, however, is by no means the only mechanism by which a network can grow, nor even the only mechanism known to generate power laws. In the remainder of this chapter we look at a number

of other models and mechanisms for the formation of networks, starting in this section with models based on node copying.

In Section 13.1 we introduced the preferential attachment mechanism and suggested a possible explanation of its origin in citation networks, that a reader perusing the literature in a given academic field would encounter references to frequently cited papers more often than references to less cited ones, and hence would be more likely to cite those frequently cited papers themselves. Another way of saying this is that, in effect, researchers are copying citations from the bibliographies of papers they read.¹⁷

Kleinberg *et al.* [269] proposed an alternative mechanism for network formation that takes this idea one step further. What if people simply copied the entire bibliography of a single paper to create the bibliography of their own paper? This would then create a new node in the network with the same pattern of outgoing edges as the node they copied from.

As we will see, this process, with slight modifications, can give rise to a power-law degree distribution. First, however, we note that the process as stated has some problems. To begin with, it's clearly rather far-fetched. Authors of papers do take note of who other authors have cited, but it seems unlikely that an author would copy the entire bibliography from someone else's paper. Moreover, if they did just copy the entire bibliography, then previously cited papers would get new citations as a result, but there would be no way for papers to receive citations if they had never been cited before.

Both of these problems can be solved by changing the model a little. Instead of assuming that the bibliography of the new paper is copied wholesale from the bibliography of an older one, let us assume that only some fraction of the entries are copied. The remainder of the new bibliography is filled out with references to other papers. These could be selected in a variety of ways, but a simple choice (and the only one we examine here) is to select them uniformly at random from the entire network.

This modification ensures that bibliographies are now no longer copied in their entirety and that papers with no previous citations have a chance of being

¹⁷We use the word "copying" figuratively here, but in fact there is evidence to suggest that some people really do just copy citations from other papers, possibly without even looking at the cited paper. Simkin and Roychowdhury [427, 428] have noted that there is a statistically surprising regularity to the typographical errors people make when citing papers. For instance, many different authors will use the same wrong page number in citing a particular paper, which suggests that, rather than copying the citation from the paper itself, they have copied it from an erroneous entry in another bibliography. This does not prove that they did not read the paper in question, but it makes it more likely—if they had actually looked up the paper, there is a good chance they would have noticed that they had the wrong page number.

cited. The model is still not a very plausible model of a real citation network, but, like Price's preferential attachment model (which is also not very realistic), it can be regarded as a simplified and tractable version of a node copying mechanism that allows us to investigate quantitatively the consequences of that mechanism.¹⁸ Moreover, the model may be quite good as a model of some other types of networks, particularly biological networks, as discussed at the end of this section.

The precise definition of the model is as follows. We will suppose for the sake of simplicity that each new node added to the network has the same out-degree c . In the language of citations, the bibliographies are all the same size. For each node added we choose uniformly at random a previous node and go one by one through the c entries in the bibliography of that previous node. For each entry we either (a) with probability γ copy that entry to the bibliography of the new node or else (b) with probability $1 - \gamma$ add a citation to a node chosen uniformly at random from the entire network. The end result is a bibliography for the new node in which, on average, γc of the entries are copied from the old node and the remainder are chosen at random. In effect, we have made an imperfect copy of the old node in which the destinations of some fraction of the outgoing edges have been randomly reassigned.

We also need to specify the starting state of the network, but, as with our preferential attachment models, it turns out that the asymptotic properties of the network do not depend on the state we choose, within reason. We could, for instance, specify a starting network consisting of some number $n_0 > c$ nodes in which each points randomly to c of the others.

We can solve for the degree distribution of the network generated by this model as follows. Let us ask what the probability is that an existing node i receives a new incoming edge upon the addition of a new node to the network. For i to receive a new edge, one of two things has to happen. Either the newly added node happens to copy connections from a node that already points to i , in which case with probability γ the connection to i will get copied, or, alternatively, i could be one of the nodes chosen at random to receive a new edge. Let us treat these two processes separately.

Suppose that our node i has in-degree q_i , meaning that q_i other nodes have

¹⁸Kleinberg *et al.* themselves proposed a different model of the copying process, but their model is quite complex and doesn't lend itself easily to analysis. The simplified model described here, which is similar to later models proposed by Solé *et al.* [432] and Vázquez *et al.* [456], embodies the important features of the process while remaining relatively tractable. We note also that Kleinberg *et al.* were not, in fact, concerned with citation networks in their paper. Their focus was on the World Wide Web. We use the language of citation networks here to emphasize the parallels with Price's model, but the discussion could equally have been framed in the language of the Web.

links to node i . The probability that a newly added node will choose to copy its links from one of these nodes is q_i/n , since the source for the copies is chosen uniformly at random from the whole network. And the chance that the link from the chosen node to i gets copied is γ , for a total probability of $\gamma q_i/n$.

At the same time, the average number of random links that a newly added node makes—ones not copied from a previous node—is $1 - \gamma$ for each of its c outgoing edges, or $(1 - \gamma)c$ in total. And the probability that our node i happens to be the target of a particular one of these random links is $1/n$, for an overall probability of $(1 - \gamma)c/n$.

Putting everything together, the total probability that node i gets a new link is¹⁹

$$\frac{\gamma q_i}{n} + \frac{(1 - \gamma)c}{n} = \frac{\gamma q_i + (1 - \gamma)c}{n}. \quad (13.112)$$

Defining $p_q(n)$ as before to be the fraction of nodes with in-degree q when the network has n nodes, the expected total number of nodes with in-degree q that receive a new edge is

$$n p_q(n) \times \frac{\gamma q + (1 - \gamma)c}{n} = [\gamma q + (1 - \gamma)c] p_q(n). \quad (13.113)$$

But now we notice an interesting fact. If we define a new constant a by

$$a = c \left(\frac{1}{\gamma} - 1 \right), \quad (13.114)$$

then

$$\gamma = \frac{c}{c + a} \quad (13.115)$$

and Eq. (13.113) becomes

$$[\gamma q + (1 - \gamma)c] p_q(n) = \frac{c(q + a)}{c + a} p_q(n), \quad (13.116)$$

which is exactly the same as the probability (13.2) for the equivalent quantity in Price's model. We can now use this probability to write down a master equation for the evolution of the degree distribution $p_q(n)$, which will be precisely the same as the master equation (13.5) for Price's model and all subsequent developments follow through just as in Section 13.1. The end result is that the degree distribution for our node copying model behaves precisely the same as

¹⁹Technically, this expression gives the expected number of new edges the node receives rather than the probability of receiving a new edge. However, in the limit of large n the two become the same.

in the Price model, but with a value of a specified now in terms of the parameter γ by Eq. (13.114). Thus, for example, the degree distribution in the limit of large n will obey Eq. (13.21) and hence will asymptotically follow a power law with exponent α given by Eq. (13.27) to be

$$\alpha = 2 + \frac{a}{c} = 1 + \frac{1}{\gamma}. \quad (13.117)$$

This gives exponents in the range from 2 to ∞ , with the value depending on how faithfully nodes are copied. Faithful copies (γ close to one) give exponents close to two, while sloppy copies give exponents that can be arbitrarily large. Other properties of Price's model carry over as well, such as the distribution of in-degree as a function of age given in Eq. (13.55).

This is not to say, however, that the node copying model generates networks identical in every respect to the preferential attachment model. With node copying, for instance, many of the links a newly appearing node makes are copied from the same preexisting node and hence most nodes in the network will have connections that are similar to those of at least one other node. In preferential attachment models, on the other hand, there is as a rule no such similarity between the connections of different nodes—each link is chosen independently from the available possibilities at the time it is created. The two networks therefore, while they may have the same degree distribution, are different in the details of their structure.

In addition to being interesting in its own right, the node copying model serves as a useful cautionary tale concerning the mechanisms of network formation. We have seen that many real networks have degree distributions that follow a power law, at least approximately, and that preferential attachment models can generate power-law degree distributions. A natural conclusion is that real networks are the product of preferential attachment processes, and this may indeed be correct. We should be careful, however, not to jump immediately to conclusions because, as we have now seen, there exists at least one other mechanism—node copying—that produces precisely the same degree distribution. Without further information we have no way of telling which of these mechanisms is the correct one, or whether some other third mechanism that we have not yet thought of is at work.

One could, in principle, examine details of the structure of specific real-world networks in an attempt to tell which, if either, of our two mechanisms is the better model for their creation. For instance, one might examine a network to see whether there appear to be pairs of nodes whose outgoing connections are approximate copies of one another. In fact, in real citation networks it turns out that there are many such pairs, an observation that appears to lend weight

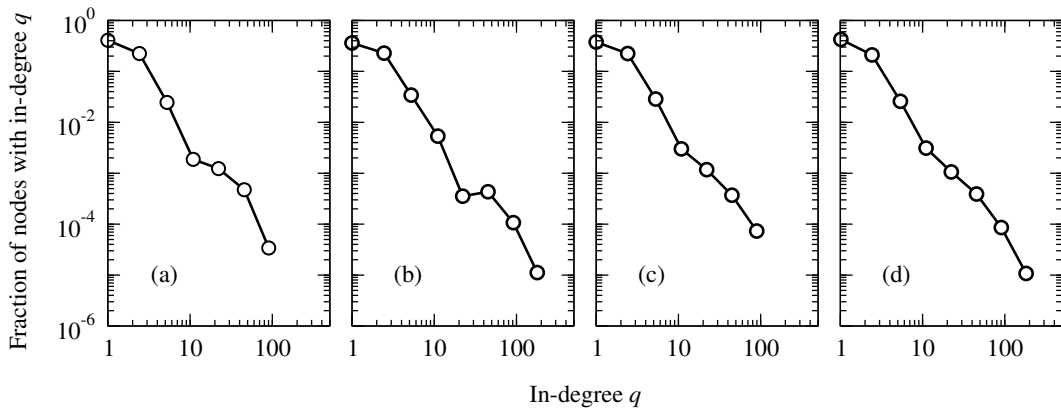


Figure 13.8: Distribution of in-degrees in the metabolic networks of various organisms. Jeong *et al.* [252] examined the degree distributions of the known portions of the metabolic networks of 43 organisms, finding some of them to follow power laws, at least approximately. Shown here are the in-degree distributions for (a) the archaeon *Archaeoglobus fulgidus*, (b) the bacterium *Escherichia coli*, (c) the nematode *Caenorhabditis elegans*, and (d) the aggregated in-degree distribution for all 43 organisms. After Jeong *et al.* [252].

to the node copying scenario. However, we must remember that both of our models are much simplified and it's likely that neither of them is an accurate representation of the way real networks are created. A simple explanation for nodes in citation networks with similar patterns of links is that they correspond to papers on similar topics and so tend to cite the same literature; there is no need to assume that one of them copied from the other. As a result, it may not be possible to distinguish firmly between preferential attachment and node copying in many cases.

There are, however, some cases where preferential attachment appears to be an implausible candidate to explain the structure of a network, and in some of these cases node copying is a promising alternative. A good example comes from the realm of biology, where node copying is considered a strong candidate to explain the structure of metabolic networks and protein–protein interaction networks. As discussed in Chapter 5, these are networks of chemical and physical interactions between molecules in the cell and, although our knowledge of their structure is currently incomplete, there is at least tentative evidence to suggest that they have power-law degree distributions—see Fig. 13.8 and Refs. [250, 252]. It seems unlikely, however, that preferential attachment is the cause of these power laws: there is no obvious mechanism by which preferential attachment would take place in this context. Node copying, on the other hand, may be a reasonable candidate.

Consider, for example, a protein interaction network. As described in Section 5.1.3, proteins in the cell are created by the processes of molecular transcription and translation from codes stored in the cell's DNA. The section of code that defines a single protein is called a gene and it turns out that genes are sometimes inadvertently copied when cells reproduce.

When a cell splits in two to reproduce, its DNA is copied so that each half of the split cell will have a complete copy. The cellular machinery responsible for the copying is highly reliable, but not perfect. Very occasionally, a section of DNA will be copied twice, giving rise to a repeated section, which can mean that the new cell has two copies of a certain gene or genes where the old cell had only one. Many examples of such repeated sections are known in the human genome and the genomes of other animals and plants.

Another type of copying error is the *point mutation*, whereby individual nucleotides—letters in the DNA code—are copied incorrectly. Over the course of many cell divisions, point mutations can accumulate, and as a result two initially identical versions of the same gene can diverge, with some fraction of their bases changed to new and (roughly speaking) random values. These processes typically happen slowly over the course of evolutionary time, taking thousands or even millions of years. The end result, however, is that a gene is copied and then mutated to be slightly different from the original.

These processes, sometimes called *duplication–divergence* processes, are reflected in the network of protein interactions. Initially, both copies of a duplicated gene generate the same protein, but the subsequent mutation of one or both copies can result in two slightly different versions of the protein, different enough in some cases to also have slightly different sets of interactions in the network. Some interactions may be common to both proteins but, just as in our node copying model, some may also be different.

This picture is made more plausible by the fact that changes in genes are not purely random but are subject to Darwinian selection under which some gene mutations are more advantageous than others. A cell with two copies of a particular gene may gain a selective advantage if those copies are slightly different, rather than needlessly duplicating functionality that a single copy alone could provide. Thus it seems possible that nature may actually favor duplicated genes that generate slightly different proteins, and hence slightly different sets of network connections. Moreover, an examination of the data for real-world protein–protein interaction networks turns up many examples of pairs of proteins that are similar but not identical in their patterns of interactions, and gene duplication is widely believed to be the cause.

Several models of node copying and mutation in biological networks have been proposed and studied. The models of Kim *et al.* [262] and Solé *et al.* [432],

for example, are both very similar to the model described in this section, the main difference being that they are models of undirected networks rather than directed ones. Another model, put forward by Vázquez *et al.* [456], is also similar but includes a mechanism whereby the connections of the *copied* node can be changed as well as those of the copying node. Although the latter mechanism would make little sense in a model of a citation network (the bibliography of a paper never changes after publication), it is appropriate in the biological context, where all genes are potentially mutating all the time.

13.6 NETWORK OPTIMIZATION MODELS

In the models we've looked at so far in this chapter, network structure is determined by the way in which the network grows—how newly added nodes connect to others, where newly added edges get placed, and so forth. Furthermore, the structure of these networks is for the most part a result of a succession of random processes, often decentralized and quite blind to the large-scale structure they are creating.

A contrasting mechanism of network formation, important in certain situations, is structural optimization. In some cases, such as transportation networks (Section 2.4) or distribution networks (Section 2.5), a network may be specifically designed to achieve a particular goal or goals, such as the delivery of mail or packages around the country or the transportation of passengers to their destinations.

Consider, for example, networks of airline routes. Such networks are typically built around a hub-and-spoke arrangement with a small number of busy airport hubs and a large number of minor destinations.²⁰ (Package delivery companies also use a similar scheme.) The reason for this arrangement is that it makes little sense to fly airplanes directly between minor destinations—there will typically be very few passengers interested in the service and the planes will be mostly empty. If all flights in and out of minor destinations are to and from major hubs, on the other hand, one concentrates the passengers on those routes, ensuring fuller planes while still giving the passengers a reasonably short journey. In other words, the hub-and-spoke design of the airline networks *optimizes* the network, making it more efficient, and hence more profitable, for the airline. Thus, in this case the structure of the network is explained

²⁰This is a relatively recent development, at least in the United States, where industry regulations made the hub-and-spoke system impractical until 1978. After regulations were lifted, the hub-and-spoke system was rapidly adopted by most of the major airlines. Hub-and-spoke systems were also adopted by the package delivery industry around the same time.

not by a growth mechanism but by the fact that the network has been designed to optimize certain characteristics. In this section we look at some models of network optimization.

13.6.1 TRADE-OFFS BETWEEN TRAVEL TIME AND COST

The example of an airline network introduced in the previous section is a good place for us to start. Airline networks are, in fact, highly optimized: the airline industry runs on very small (sometimes even negative) profit margins, and optimization of operations to trim even a tiny percentage off their enormous costs can make a substantial difference to the bottom line. Airlines employ large staffs of researchers whose sole task is to find new ways to optimize aspects of their business, including particularly their network of routes.

At the same time, airlines need to keep their customers happy, if they are to avoid losing market share to their competitors. This means, for instance, that they need to provide short, quick routes between as many pairs of destinations as possible—travelers are strongly averse to long journeys that wear them out or waste their time. The twin goals of cost-efficient operation and short routes are to some extent at odds with one another. The quickest way to get passengers from any place to any other, for example, would be to fly separate planes between every pair of airports, but this would be immensely costly. The observed structure of real airline networks is a compromise response to the conflicting needs of the company and its passengers.

The optimization problems faced by real airlines are, inevitably, hugely complex, involving as they do organizations with thousands of employees, billions of dollars worth of material resources, and rapidly changing parameters such as fuel costs, consumer demand, and the nature of the competition. Nonetheless, there is insight to be gained by creating and studying simplified models of the optimization process in the same way that simple models of, for example, citation networks can grant us insight despite the many features of real citation processes that they omit.

One of the simplest models of network optimization is that proposed by Ferrer i Cancho and Solé [176], which balances two elements of exactly the types discussed above. In this model the cost of maintaining and operating the network is represented by the number of edges m in the network. This would be equivalent to saying that the cost of running an airline is proportional to the number of routes it operates. Obviously this is a simplification of the real situation, but let us accept it for the moment and see where it leads. The customer satisfaction half of the equation is represented by the mean network distance ℓ between node pairs. In our airline example ℓ would be the average

number of legs required to journey from one point to another, which is certainly one element of customer satisfaction, though not the only one. Technically, ℓ is a *dissatisfaction* measure, since large values correspond to disgruntled customers.

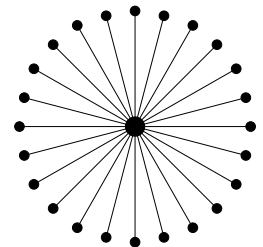
We would like to design a network that minimizes both m and ℓ but this is, in general, not possible: the minimum value of ℓ is achieved by placing an edge between every pair of nodes, but this maximizes the value of m . Thus our two goals are, as discussed above, at odds with one another and the best we can hope for is a reasonable compromise between them. In search of such a compromise, Ferrer i Cancho and Solé studied the quality function

$$E(m, \ell) = \lambda m + (1 - \lambda)\ell, \quad (13.118)$$

where λ is a parameter in the range $0 \leq \lambda \leq 1$. For any given network and a given value of λ we can calculate $E(m, \ell)$. The value of ℓ , for instance, can be computed using the “all-pairs” variant of the breadth-first search algorithm described in Section 8.5.4. Ferrer i Cancho and Solé considered networks of a given number of nodes n and then asked what happens when we try to minimize $E(m, \ell)$ by varying the position of the edges in that network to find the smallest value possible. If $\lambda = 1$, then $E = m$ and this process is equivalent to just minimizing the number of edges without regard for path lengths. If $\lambda = 0$ then $E = \ell$ and we are minimizing only average path length without regard for m . For values in between, we are striking a balance between number of edges and path length, with the precise weight of each term controlled by our choice of λ .

At some level, this model is a trivial one. If we adopt the common convention that the distance between two nodes is infinite if there is no path between them—i.e., if they lie in different components—then the minimum value of E must occur for a connected network, a network with just one component. Observe also that the minimum value of m for a connected network is $m = n - 1$, where n is the number of nodes. This is the value for a tree, which is the connected network with the smallest number of edges (see Section 6.8).

Provided λ is reasonably large, so that we place a moderate amount of weight on minimizing m , the network with the best value of $E(m, \ell)$ is then found by giving m its minimum value of $n - 1$, which constrains the network to be a tree, and searching through the set of possible trees to find the one that minimizes ℓ . But the latter task has a simple, known solution: the minimum value of ℓ among trees with n nodes is obtained by the *star graph*, the network in which there is a single central hub connected by edges to each of the $n - 1$ remaining nodes. To see this, note that by definition there are in any network exactly m pairs of nodes that are separated by distance one—the pairs that are directly connected by an edge—which means that in a tree there are $n - 1$ such



A star graph of 25 nodes.

pairs. Among the set of all trees of a given size, therefore, the value of the mean distance ℓ is governed by the numbers of pairs separated by distance two or more, since the number with distance one is fixed. But in the star graph all other pairs have distance exactly two—the shortest (indeed only) path from any (non-hub) node to any other is the path of length two via the hub. Thus there can be no other tree with a smaller value of ℓ .

Thus, for sufficiently large λ , the network with the smallest value of the quality function (13.118) is always the star graph. This is satisfying to some extent. It offers a simple explanation for why the hub-and-spoke system is so efficient: it offers short path lengths while still being economical in terms of number of edges. But it is also, as we have said, somewhat trivial. The model shows essentially only the one behavior. For smaller values of λ other behaviors are possible, but it turns out that the value of λ must be *really* small: non-star-graph solutions only appear when²¹

$$\lambda < \frac{2}{n^2 + 2}. \tag{13.119}$$

Since the expression on the right-hand side dwindles rapidly as n becomes large, the optimal network is a star graph for almost all values of λ , even for networks of quite modest size.

In their paper, however, Ferrer i Cancho and Solé did not perform precisely the calculation we have done here. Instead, they took a different and interesting approach, in which they looked for *local minima* of $E(m, \ell)$, rather than the global minimum. They did this numerically, starting with a random network,

²¹The derivation of this result is as follows. If λ is sufficiently large then, as we have shown, the optimal network is the star graph. If we now reduce λ slowly, then at some point we enter a regime in which the cost of adding an edge between the “spoke” nodes of the star graph is sufficiently offset by the corresponding reduction in the mean distance that it becomes worthwhile to add such edges. To calculate the point at which this happens let us take our star graph and add to it some number r of extra edges. Necessarily these edges fall between the spoke nodes, since there is nowhere else for them to fall, and in doing so they form paths of length one between pairs of nodes whose previous shortest connecting path was of length two. The shortest paths between no other nodes are affected by the addition. Thus the total number of node pairs connected by paths of length 1 is $n - 1 + r$ and all the rest have paths of length two. Then the mean distance, as defined in Eq. (10.2), is

$$\ell = \frac{1}{n^2} \sum_{ij} d_{ij} = 2 \frac{(n - 1 + r) + 2[\frac{1}{2}n(n - 1) - (n - 1 + r)]}{n^2} = 2 \frac{(n - 1)^2 - r}{n^2}.$$

(The leading factor of two comes from the fact that the sum over i, j counts each pair of nodes twice.) Substituting this expression, along with $m = n - 1 + r$, into Eq. (13.118) then gives

$$E = \lambda(n - 1 + r) + 2(1 - \lambda) \frac{(n - 1)^2 - r}{n^2} = \text{constant} + \left[\lambda + \frac{2(\lambda - 1)}{n^2} \right] r.$$

repeatedly choosing a pair of nodes at random, and either connecting them by an edge if they were not already connected or deleting the edge between them if they were. Then they compared the value of E before and after the change. If E decreased or stayed the same, they kept the change. If not, they reverted back to the state of the network before the change. The whole procedure was then repeated until the value of E stopped improving, meaning in practice that a long string of attempted changes were rejected because they increased E .

An algorithm of this kind is called a *random hill climber* or *greedy algorithm*. The networks it finds are networks for which the value of E cannot be reduced any further by the addition or removal of any single edge. This does not mean that networks with lower values of E don't exist: there may be states of the network that differ by more than one edge—the addition and deletion of whole regions of the network—that have better values of E . But if so, the algorithm will not find them. It comes to a halt at a local minimum where no single-edge change will improve the value of E .

When studied in this way, the model shows an interesting behavior. For large values of λ , where the addition of an edge costs a great deal in terms of the value of E , the algorithm rapidly runs into trouble and cannot find a way to improve the network, long before it gets anywhere close to the optimum hub-and-spoke arrangement. When λ is small, on the other hand, the algorithm typically manages to find the star graph solution or something like it. The result is a spectrum of networks that range from a random-looking tree to something close to a star graph, as shown in Fig. 13.9.

What's more, Ferrer i Cancho and Solé found that the degree distributions of their networks show interesting behavior, passing from an exponential distribution for large λ , through a transition point with a power-law degree distribution, to approximately star-like graphs for small λ in which one node gets a finite fraction of all the edges and the remaining nodes have low degree. This spectrum is reminiscent of the behavior of continuous phase transitions such as the transition at which a giant component appears in a random graph (see Section 11.5). Recall that on one side of that transition there is an exponential

This will decrease with growing r only if the quantity in square brackets [...] is negative, i.e., if

$$\lambda < \frac{2}{n^2 + 2}.$$

If this condition is satisfied then it becomes advantageous to add edges between the spoke nodes, and to keep on doing so until the network becomes a complete graph, with every node connected to every other. Thus there is a discontinuous transition between two behaviors—the star graph and the complete graph—at the point $\lambda = 2/(n^2 + 2)$. Real distribution and transportation networks appear to be in the star graph regime.

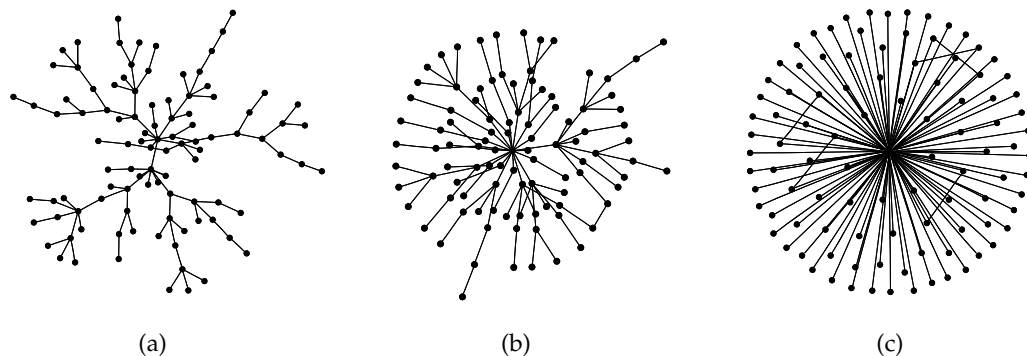


Figure 13.9: Networks generated by the optimization model of Ferrer i Cancho and Solé. The optimization model described in the text generates a range of networks, all trees or approximate trees, running from (a) distributed networks with exponential degree distributions, through (b) power-law degree distributions, to (c) star-like graphs in which there is just one major hub. Figure adapted from [176]. Original figure Copyright 2003 Springer-Verlag Berlin Heidelberg. Reproduced with permission of Springer Nature.

distribution of component sizes, while on the other side one component gets a finite fraction of all nodes and the rest are small.

Sadly, this observation does not go any further than an intriguing hint. The work of Ferrer i Cancho and Solé is entirely numerical and no analytic treatment of the model has been given to date. In addition, there are some other problems with the model. In particular, it is not clear why one should look at local minima of E rather than global ones: the researchers who work for real airlines are certainly capable of realizing when they are stuck in a local minimum and better profits are available by changing the network in some substantial way that moves them to a different and better state. It seems likely therefore that, to the extent that real networks show interesting structural behavior of the type observed here, it is not a result of getting stuck in local minima and hence that a model with a different approach is needed.

One such model, proposed by Gastner and Newman [202], generalizes that of Ferrer i Cancho and Solé by considering not only the number of legs in a journey but also the geographic distance traveled. Suppose that airline travelers are principally concerned not with the number of legs in their journey but with the total time it takes them to travel from origin to destination. Number of legs can be regarded as a simple proxy for travel time, but a better proxy would be to take the length of those legs into account as well as their number. The travel time contributed to a journey by one leg is composed of the time spent in the airport (checking in, waiting, embarking, taxiing, disembarking, etc.) plus the

time spent in the air. A simple formula would be to assume that the former is roughly constant, regardless of the distance being traveled, while the latter is roughly proportional to the distance traveled. Thus, the time taken by a leg from node i to node j in the network would be

$$t_{ij} = \mu + \nu r_{ij}, \quad (13.120)$$

where μ and ν are constants and r_{ij} is the distance flown from i to j . By varying the values of μ and ν , we can place more or less emphasis on the fixed airport time cost versus the time spent in the air.

Gastner and Newman used this expression for travel time in place of the simple hop-count of the model of Ferrer i Cancho and Solé, redefining ℓ to be the average shortest-path distance between pairs of nodes when distances are measured in terms of travel time. The quality function E is then defined as before, Eq. (13.118), but using this new definition of ℓ .

Despite the superficial similarity between this model and that of Ferrer i Cancho and Solé, there is a crucial difference between the two: the model of Gastner and Newman depends on actual spatial distances between airports and hence requires that the nodes of the network be placed at some set of positions on a map. The model of Ferrer i Cancho and Solé by contrast depends only on the network topology and has no spatial element. There are various ways in which the nodes can be positioned on the map. Gastner and Newman specifically considered the map of the United States and took the real US population distribution into account, placing nodes with greater density in areas with greater populations. However, while this adds a level of realism to the calculation, the interesting behavior of their model can be seen without going so far. In the examples given here we consider a fictional map in which nodes are just placed uniformly at random in a square with periodic boundary conditions.

Another important difference between the two models is that Gastner and Newman considered the global optimum of the quality function rather than local optima as Ferrer i Cancho and Solé did. The global optimum is hard to find, however, so usually we must make do with approximate optima calculated numerically. Gastner and Newman were able to find good approximations to the global optimum using the numerical optimization technique known as simulated annealing, but we should bear in mind that they are only approximations.

Figure 13.10 shows optimal or approximately optimal networks for various values of the parameters μ and ν . The leftmost frames of the figure correspond to small μ and large ν , meaning that the cost to the traveler of a trip is roughly proportional to the total mileage traveled and the number of legs has little effect. In this case, the best networks are ones that allow travelers to travel in

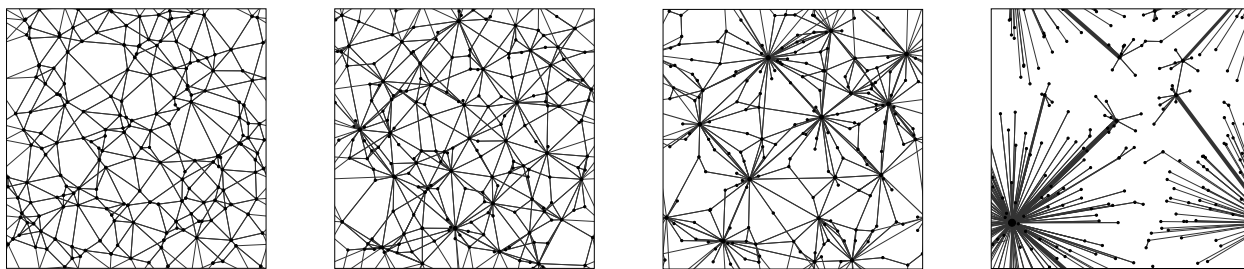


Figure 13.10: Networks generated by the spatial network model of Gastner and Newman. The four frames show networks that optimize or nearly optimize the quality function, Eq. (13.118), with ℓ defined according to the prescription of Gastner and Newman [202] in which distances are replaced by the approximate travel time required to traverse edges. Travel time has two components, a fixed cost per edge and a cost that increases with the Euclidean length of the edge. The frames show the resulting networks as the relative weight of these two components is varied between the extremes represented by the network on the left, for which all of the weight is on the Euclidean length, and the network on the right, for which cost is the same for all edges. The resulting structures range from road-like in the former case, to airline-like in the latter. Adapted from Gastner [201].

roughly straight lines from any origin to any destination. As the figure shows, the networks are roughly planar in appearance. They look reminiscent of road networks, rather than airline routes, and this is no coincidence. Travel times for road travelers are indeed dominated by total mileage: there is almost no “per leg” cost associated with road travel, since it takes only a few seconds to turn from one road onto another. It is satisfying to see therefore that the simple model of Gastner and Newman generates networks that look rather like real road maps in this limit.

The rightmost frames in the figure show optimal networks for large μ and small ν —the case where it is mostly the number of legs that matters and the length of those legs is relatively unimportant. As we saw also for the model of Ferrer i Cancho and Solé, the best networks in this case are star-like hub-and-spoke networks.

Thus the model interpolates between road-like and airline-like networks as the parameters are varied from one extreme to the other. Note that the parameter λ governing the cost of building or maintaining the network is held constant in Fig. 13.10. In principle, we could vary this parameter too, which would affect the total number of edges in the network. For higher λ sparser networks with fewer edges would be favored, while for lower λ we would see denser networks. The work of Gastner and Newman still suffers from the drawback that the results are numerical only. Some analytic results for the model have, however, been derived by Aldous [18].

EXERCISES

13.1 Consider the growing network model of Price introduced in Section 13.1.

- Using the results given in this chapter write down an expression in terms of the parameters a and c for the expected in-degree of the i th node added to the network just before the j th node is added, where $i < j$. You can assume j is large.
- Hence show that the average probability of a directed edge from j to i in a network with n nodes, where $n \geq j$, is

$$P_{ij} = \frac{ca}{c+a} i^{-c/(c+a)} (j-1)^{-a/(c+a)}.$$

13.2 Within a set of papers on a particular topic it is found that the average paper cites 30 others in the set. Moreover the network of citations among the papers appears to be scale-free, with power-law exponent $\alpha = 3$. Let us assume that the network is well described by the preferential attachment model of Price, discussed in Section 13.1.

- What is the average number of citations received by a paper?
- On average what fraction of papers receive no citations at all?
- On average what fraction of papers receive 100 or more citations?
- The set consists of 10 000 papers. What is the probability that the 100th paper published has no citations? What is the probability that the 100th-to-last paper published has no citations?

13.3 Consider Price's model as a model of a citation network, applied to publications in a single field, a field that is currently, say, five years old.

- Suppose that you are the author of the tenth paper published in the field. Assuming papers in the field are published at a constant rate, how long will it be from now before the expected number of citations your paper has within the field is equal to the expected number that the first paper published currently has?
- Derive an expression for the average number of citations received by a paper published between times τ_1 and τ_2 , where time is defined as in Eq. (13.42).
- Reasonable values of the model parameters for real citation networks are $c = 20$ and $a = 5$. For these parameter choices, what is the average number of citations to a paper in the first 10% of those published? And what is the average number for a paper in the last 10%?

These perhaps surprising numbers are examples of the first mover advantage discussed in Section 13.3—the substantial bias of citation numbers in favor of the first papers published in a field.

13.4 Recall the master equations (13.7) and (13.8) for Price's model in the limit of large n :

$$p_q = \frac{c}{c+a} [(q-1+a)p_{q-1} - (q+a)p_q] \quad \text{for } q \geq 1,$$

$$p_0 = 1 - \frac{ca}{c+a} p_0 \quad \text{for } q = 0.$$

- a) Write down the special case of these equations for $c = a = 1$.
 b) Show that the in-degree generating function $g_0(z) = \sum_{q=0}^{\infty} p_q z^q$ for this case satisfies the differential equation

$$g_0(z) = 1 + \frac{1}{2}(z-1)[z g_0'(z) + g_0(z)].$$

- c) Show that the function

$$h(z) = \frac{z^3 g_0(z)}{(1-z)^2}$$

satisfies

$$\frac{dh}{dz} = \frac{2z^2}{(1-z)^3}.$$

- d) Hence find a closed-form solution for the generating function $g_0(z)$. Confirm that your solution has the correct limiting values $g_0(0) = p_0$ and $g_0(1) = 1$.
 e) Thus find a value for the mean in-degree of a node in Price's model. Is this what you expected?

13.5 Consider the following variant of the Barabási–Albert model. Nodes are added one by one to a growing undirected network, each node having initial degree c . The c edges emanating from a newly added node connect to previously existing nodes i with probability proportional to $k_i + a$, where k_i is i 's (undirected) degree and a is a constant.

- a) Given that c edges are added to the network with each node, what is the mean degree of a node in the network in the limit of large network size?
 b) Derive the master equation that gives the fraction of nodes p_k having degree k in the limit of large network size. If necessary, give an additional rate equation to cover any special-case value of k .
 c) Show that the fraction of nodes with degree c in the limit of large network size is

$$p_c = \frac{2c + a}{2c + a + c(c + a)}.$$

13.6 Consider a model of a growing directed network similar to Price's model described in Section 13.1, but without preferential attachment. That is, nodes are added one by one to the growing network and each has c outgoing edges, but those edges now attach to existing nodes chosen uniformly at random, without regard for degrees or any other node properties.

- a) Derive master equations, the equivalent of Eqs. (13.7) and (13.8), that govern the distribution of in-degrees q in the limit of large network size.
 b) Hence show that in the limit of large size the in-degrees have an exponential distribution $p_q = C r^q$, where C is a normalization constant and $r = c/(c + 1)$.

13.7 Consider a model network similar to the model of Barabási and Albert described in Section 13.2, in which undirected edges are added between nodes according to a preferential attachment rule, but suppose now that the network does not grow—it starts off with a given number n of nodes and neither gains nor loses any nodes thereafter. In this

model, starting with an initial network of n nodes and some specified arrangement of edges, we add at each step one undirected edge between two nodes, both of which are chosen at random in direct proportion to degree k . Let $p_k(m)$ be the fraction of nodes with degree k when the network has m edges.

- Show that when the network has m edges, the probability that node i will get a new edge upon the addition of the next edge is k_i/m .
- Write down a master equation giving $p_k(m+1)$ in terms of $p_{k-1}(m)$ and $p_k(m)$. Be sure to give the equation for the special case of $k=0$ also.
- Eliminate m from the master equation in favor of the mean degree $c = 2m/n$ and take the limit $n \rightarrow \infty$ with c held constant to show that $p_k(c)$ satisfies the differential equation

$$c \frac{dp_k}{dc} = (k-1)p_{k-1} - kp_k.$$

- Define a generating function $g(c, z) = \sum_{k=0}^{\infty} p_k(c) z^k$ and show that it satisfies the partial differential equation

$$c \frac{\partial g}{\partial c} + z(1-z) \frac{\partial g}{\partial z} = 0.$$

- Show that $g(c, z) = f(c - cz)$ is a solution of this differential equation, where $f(x)$ is any differentiable function of x .
- The particular choice of f depends on the initial conditions on the network. Suppose the network starts off in a state where every node has degree one, which means $c = 1$ and $g(1, z) = z$. Find the function f that corresponds to this initial condition and hence find $g(c, z)$ for all values of c and z .
- Show that, for this solution, the degree distribution as a function of c takes the form

$$p_k(c) = \frac{(c-1)^{k-1}}{c^k},$$

except for $k=0$, for which $p_0(c) = 0$ for all c .

Note that this degree distribution decays exponentially in k , implying that preferential attachment does not, in general, generate a power-law degree distribution if the network is not also growing.

13.8 Consider a model of a growing network similar to Price's model described in Section 13.1, but in which the parameter a , which governs the rate at which nodes receive new incoming links when their current in-degree is zero, varies from node to node. That is, at each step a new node with c outgoing edges is added to the network, and the probability of one of the new edges attaching to node i is proportional to $q_i + a_i$ where q_i is the current in-degree of node i and a_i is a parameter defined separately for each node but constant in time. In the context of citation networks, for example, a_i could be considered a measure of the intrinsic merit of a paper, controlling the rate at which the paper gets citations immediately after first publication (when $q_i = 0$).

- Suppose that a_i for each node is drawn at random from a stationary distribution $\rho(a)$, so that the probability of falling in the range from a to $a + da$ is $\rho(a) da$.

You can assume that the distribution has a well-defined mean. Show that, in the limit of large n , the probability that the $(n + 1)$ th node added to the network attaches to a previously added node i with in-degree q_i is $c(q_i + a_i)/n(c + \bar{a})$, where $\bar{a} = \int a \rho(a) da$ is the average value of a_i .

- b) Let $p_q(a) da$ be the probability that a node has degree q and parameter a in the range a to $a + da$ in the limit of large network size. Derive the master equations that govern $p_q(a)$ in the limit of large n . Pay special attention to the equation for $p_0(a)$, which is different from other values of q .
- c) Solve the master equations to show that

$$p_q(a) = \frac{B(q + a, 2 + \bar{a}/c)}{B(a, 1 + \bar{a}/c)} \rho(a),$$

where $B(x, y)$ is Euler's beta function. (Note the distinction between a and \bar{a} in this formula.) This distribution has a power-law tail. What is the exponent of the power law?

13.9 Consider the following simple model of a growing network. Nodes are added to a network at a rate of one per unit time. Edges are added at a mean rate of β per unit time, where β can be anywhere between zero and ∞ . (That is, in any small interval δt of time, the probability of an edge being added is $\beta \delta t$.) Edges are placed uniformly at random between any pair of nodes that exist at that time. They are never moved after they are first placed.

We are interested in the component structure of this model, which we will tackle using a master equation method. Let $a_k(n)$ be the fraction of nodes that belong to components of size k when there are n nodes in the network. Equivalently, if we choose a node at random from the n nodes currently in the network, $a_k(n)$ is the probability the node will belong to a component of size k .

- a) What is the probability that a newly appearing edge will fall between a component of size r and another of size s ? (You can assume that n is large and the probability of both ends of an edge falling in the same component is small.) Hence, what is the probability that a newly appearing edge will join together two pre-existing components to form a new one of size k ?
- b) What is the probability that a newly appearing edge joins a component of size k to a component of *any* other size, thereby creating a new component of size larger than k ?
- c) Thus write down a master equation that gives the fraction of nodes $a_k(n + 1)$ in components of size k when there are $n + 1$ nodes in total.
- d) The only exception to the previous result is that components of size 1 appear at a rate of one per unit time. Write a separate master equation for $a_1(n + 1)$.
- e) If a steady-state solution exists for the component size distribution, show that it satisfies the equations

$$(1 + 2\beta)a_1 = 1, \quad (1 + 2\beta k)a_k = \beta k \sum_{j=1}^{k-1} a_j a_{k-j}.$$

f) Multiply by z^k and sum over k from 1 to ∞ and hence show that the generating function $g(z) = \sum_k a_k z^k$ satisfies the ordinary differential equation

$$2\beta \frac{dg}{dz} = \frac{1 - g/z}{1 - g}.$$

Unfortunately, the solution to this equation is not known, so for the moment at least we do not have a complete solution for the component sizes in this model.

13.10 Check the solution given in Eq. (13.55) in two different ways:

- a) Substitute (13.55) into (13.47) and confirm directly that the former is indeed a solution of the latter for general q .
- b) Integrate Eq. (13.55) over τ from 0 to 1 and confirm that you recover the correct form for the degree distribution of Price's model, Eq. (13.21). You will probably need to use the integral formula for the Euler beta function given in Eq. (13.33).

PART IV
APPLICATIONS

CHAPTER 14

COMMUNITY STRUCTURE

A discussion of methods for identifying groups or communities of nodes within networks

THE ULTIMATE goal in studying networks is to better understand the behavior of the systems they represent. For instance, we study social networks to better understand the nature of social interactions and their implications for human experience, commerce, the spread of disease, and the structure of society. We study the Internet to better understand the flow of data traffic or why communications protocols function the way they do or how we might change the network to make it perform better. We study biochemical networks like metabolic networks because we hope it will lead to a better understanding of the complex chemical processes taking place in the cell and perhaps even to new therapies for disease or injury.

The techniques discussed in the previous chapters of this book, which focus on measurements and models of network structure, provide a solid foundation for this kind of understanding, but they are only a beginning. Our task now is to apply what we have learned to gain deeper insight into the *function* of networked systems. Unfortunately, research progress in this direction has been slower than on measurements and models, but there are some areas in which substantial advances have been made, including the development of new analysis techniques such as community detection, studies of network failure and resilience, and studies of epidemic and other spreading processes. The remaining chapters of this book are devoted to a description of our current understanding of some of these issues. We begin in this chapter by looking at

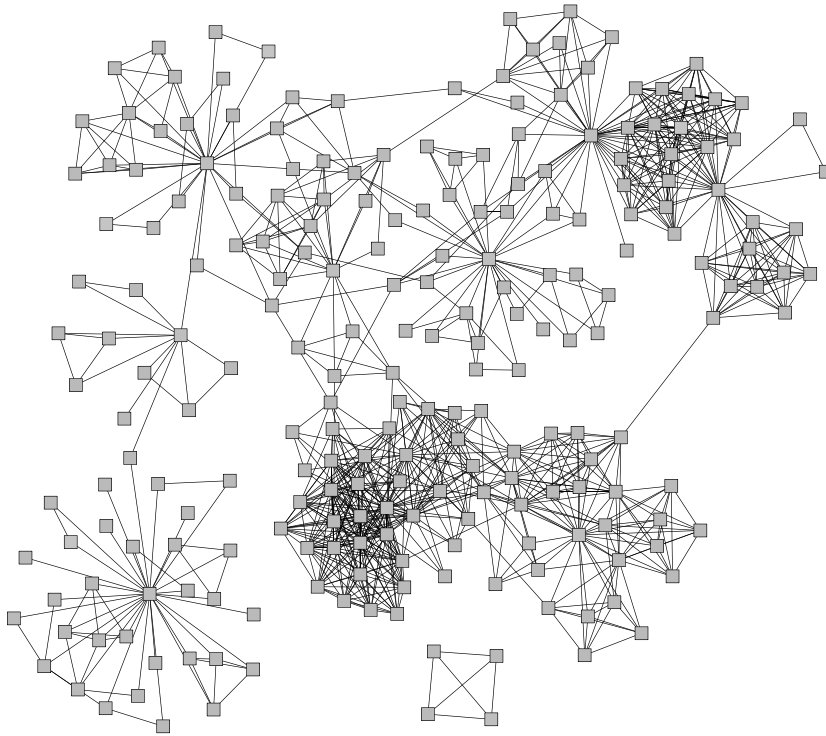


Figure 14.1: Network of coauthorships in a university department. The nodes in this network represent scientists in a university department, and edges link pairs of scientists who have coauthored scientific papers. The network has clear clusters or “community structure,” presumably reflecting divisions of interests and research groups within the department.

community detection, one of the most active areas of current networks research and one that will introduce us to a range of important approaches and phenomena, including information theory, optimization methods, and methods of statistical inference.

14.1 DIVIDING NETWORKS INTO GROUPS

Consider Fig. 14.1, which shows patterns of collaborations between scientists in a university department. Each node in this network represents a scientist and edges between nodes indicate pairs of scientists who have coauthored one or more papers together. As we can see from the figure, this network

contains a number of densely connected clusters of nodes, corresponding to groups of scientists who have worked closely together. Readers familiar with the organization of university departments in the sciences will not be surprised to learn that these clusters correspond, at least approximately, to established research groups within the department.

But suppose one did not know how university departments operate and wished to study them. By constructing a network like that of Fig. 14.1 and then observing its clustered structure, one could deduce the existence of groups within the larger department, and by further investigation could probably quickly work out how the department was organized. Thus the ability to find groups or clusters in a network can be a useful tool for revealing structure and organization within networks at a scale larger than that of a single node or a few nodes. Figure 7.12 on page 202 shows another example, in a network of friendships between US high school students. In this case the network splits into two clear groups or communities, which, as described in Section 7.7, are primarily dictated by students' ethnicity, and this structure can give us clues about the nature of the social interactions within the population represented by the network.

The occurrence of groups or communities is not limited to social networks. Clusters of nodes in a web network, for instance, might indicate groups of related web pages. Clusters of nodes in a metabolic network might indicate functional units within the network.

The ability to find groups also has another practical application: it allows us to take a large network and break it apart into smaller subsets that can be studied separately. The network in Fig. 14.1 is quite small, but others may be much larger, millions of nodes or more, making their analysis and interpretation challenging. Breaking such networks into their component clusters is a useful technique for reducing them to a manageable size. One example of this approach is in network visualization. A network with a million or more nodes can rarely be visualized in its entirety, even with the aid of good visualization software. Such networks are simply too big to be represented usefully on the screen or on paper. If the nodes of the network divide naturally into groups, however, then we can make a simpler but still useful picture by representing each group as a single node and the connections between groups as edges. An example is shown in Fig. 14.2. This simplified representation allows us to see the large-scale structure of the network without getting bogged down in the details of the individual nodes. If one wanted to see the individual nodes, one could then "zoom in" to a single group and look at its internal makeup.

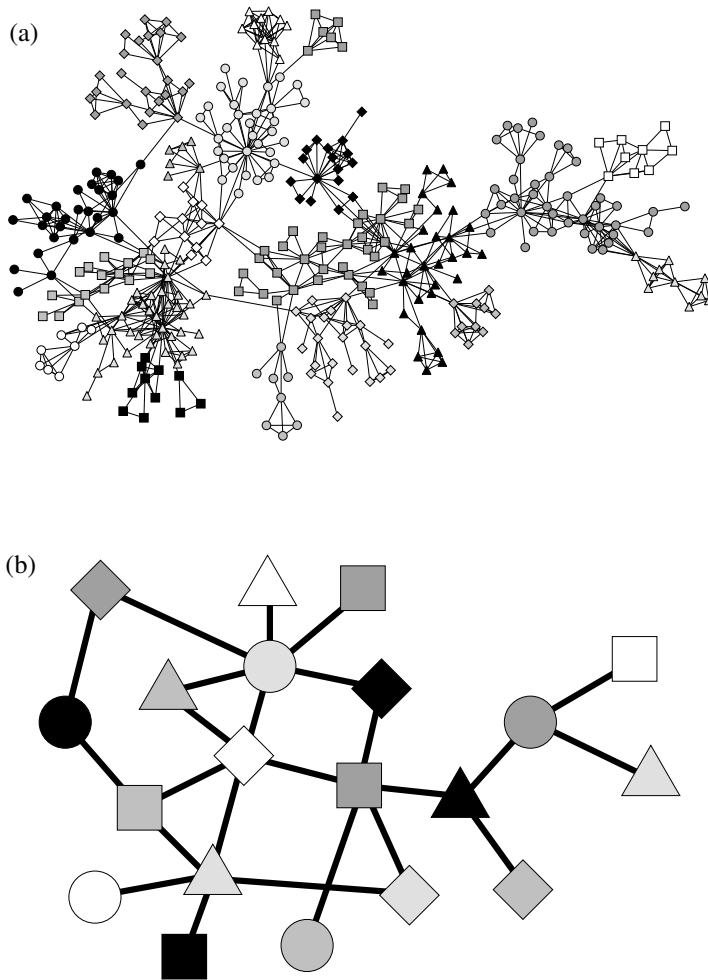


Figure 14.2: Visualization of network structure using community detection. The network in (a) is decomposed into its constituent communities, represented by the different shapes and colors. (The communities were found using the spectral modularity maximization method of Section 14.2.3.) In (b) each community in the network is represented by a single large node, with edges indicating which communities are connected to which. The overall pattern of connections in this coarse-grained representation of the network is now easy to see.

Community detection is sometimes also called “clustering,” although we largely avoid this term to prevent confusion with the other use of the word clustering introduced in Section 7.3.

A number of network visualization packages, such as Gephi,¹ have the ability to perform this kind of two-level visualization.

The problem of finding groups of nodes in networks is called *community detection*. Simple though it is to describe, community detection turns out to be a challenging task, but a number of methods have been developed that return good results in practical situations. In this chapter we describe some of the most widely used.

14.2 MODULARITY MAXIMIZATION

Part of what makes community detection difficult is that the problem is not very well posed. Loosely speaking, the goal of community detection is to find the natural divisions of a network into groups of nodes such that there are many edges within groups and few edges between them. This description, however, is vague and open to interpretation—what exactly do we mean by “many” edges or “few”? To turn community detection into a problem we can tackle quantitatively, we need to put some numbers on these concepts. There are a variety of ways to do this. We begin in this section with the most widely used approach, the method of modularity maximization.

We approach the task of community detection as an optimization problem. Consider a simple undirected, unweighted network. To every possible division of this network into communities we assign a score which is high if that division is “good,” in the above sense of having many edges within communities and few edges between them, and low if it is “bad.” Then we search through the divisions to find the one with the highest score, which we take to be the best division of the network. The success of this approach clearly relies on finding a satisfactory definition of the score that correctly gives high scores to good divisions. But this is a problem we have encountered before. In Section 7.7 we considered the phenomenon of assortative mixing in networks, in which nodes with similar characteristics tend to be connected by edges. There we introduced the measure of assortative mixing known as modularity, which has a high value when connections are primarily between nodes of the same type and a low value when they are not. This is precisely the kind of measure we need for our community detection problem. We can think of the members of our groups as different types of nodes, and good divisions of the network into communities are precisely those that have high values of the corresponding modularity. Thus, one way to detect communities in networks is to look for

¹See <https://gephi.org>

the division that has the highest modularity score. This is the method of modularity maximization.

Let us initially consider the simplest possible community detection problem, the problem of dividing a given network into just two groups or communities. That is, we will for the moment assume that we know there to be exactly two groups in the network and our only task is to determine which nodes belong to which group. The groups can be of any size we like, but every node must belong to one group or the other, so the sizes of the two groups sum to n , the size of the network as a whole.

Applying the modularity maximization method directly to this problem entails going through all possible assignments of the nodes to the two groups and finding the one with the highest modularity. Unfortunately, however, this turns out to be a difficult task. There are 2^n ways to divide n nodes into two groups and this number becomes large very quickly as n grows.² Even for a relatively small network of just a hundred nodes, for example, there are $2^{100} \approx 10^{30}$ possible divisions, and calculating the modularity for all of them would be well beyond the capabilities of any current computer.

One might wonder whether it is possible to get around this problem by clever programming. Brute-force enumeration of all possible divisions is not a very imaginative way to find the best one. Could one perhaps find some way to limit one's search to only those divisions that have a chance of being the best? Unfortunately, it is believed that this is not possible: there are fundamental results from computer science that tell us that no such algorithm will ever be able to find the best division of the network in all cases. Either an algorithm can be clever and run quickly, but will fail to find the optimal answer in some (and perhaps most) cases, or it always finds the optimal answer but takes an impractical length of time to do it. These are the only options.³

This is not to say, however, that clever algorithms for modularity maximiza-

²Really there are only 2^{n-1} divisions: one puts the first node in either group and then one is left to decide for each of the $n - 1$ remaining nodes whether they go in the same group as the first one, or the other group. The argument still holds however: 2^{n-1} grows quickly with n and the number of divisions will rapidly outpace our ability to search through them all, even for relatively modest values of n .

³Technically, this statement has not actually been proved. Its truth hinges on the assumption that two fundamental classes of computational problems, called P and NP, are not the same. Although this assumption is universally believed to be true—our world would pretty much fall apart if it weren't—no one has yet proved it, nor even has any idea about where to start. Readers interested in the fascinating branch of computer science that deals with problems of this kind are encouraged to look, for example, at the book by Moore and Mertens [340]. For applications to the specific problem of modularity maximization, see Brandes *et al.* [81].

tion do not exist or that they don't give useful answers. Even algorithms that fail to find the very best division of a network may still find a pretty good one, and for many practical purposes pretty good is good enough. The goal of essentially all practical algorithms is just to find a "pretty good" division in this sense. Algorithms that find approximate, but acceptable, solutions to problems in this way are called *heuristic algorithms* or just *heuristics*. A range of heuristic algorithms have been tried for the modularity maximization problem. We will discuss some of them shortly.

14.2.1 THE FORM OF THE MODULARITY FUNCTION

In Section 7.7.1 we wrote an expression, Eq. (7.54), for the modularity of a division of a network thus:

$$Q = \frac{1}{2m} \sum_{ij} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{g_i g_j} = \frac{1}{2m} \sum_{ij} B_{ij} \delta_{g_i g_j}, \quad (14.1)$$

where we have numbered the groups or communities (for instance starting from 1) and g_i is the number of the group to which node i belongs, δ_{ij} is the Kronecker delta, and

$$B_{ij} = A_{ij} - \frac{k_i k_j}{2m}. \quad (14.2)$$

Note that B_{ij} has the property that its sum with respect to either of its indices is zero:

$$\sum_i B_{ij} = \sum_i A_{ij} - \frac{k_j}{2m} \sum_i k_i = k_j - \frac{k_j}{2m} 2m = 0, \quad (14.3)$$

$$\sum_j B_{ij} = \sum_j A_{ij} - \frac{k_i}{2m} \sum_j k_j = k_i - \frac{k_i}{2m} 2m = 0, \quad (14.4)$$

where we have made use of Eq. (6.13). This property will be important shortly.

Equation (14.1) is completely general, applying to divisions of a network into any number of groups, but let us consider again the case of just two groups, for which we can usefully rewrite the equation by defining the new quantities⁴

$$s_i = \begin{cases} +1 & \text{if node } i \text{ belongs to group 1,} \\ -1 & \text{if node } i \text{ belongs to group 2.} \end{cases} \quad (14.5)$$

⁴To a physicist, the quantities s_i are "Ising spins" on the nodes of the network and Eq. (14.7) takes the form of the Hamiltonian of an Ising spin-glass. It was in part this connection that inspired the formulation presented here, but it is not necessary to be familiar with the physics to understand how modularity maximization works.

With this definition, the quantity $\frac{1}{2}(s_i s_j + 1)$ is 1 if i and j are in the same group and zero otherwise, so that

$$\delta_{g_i, g_j} = \frac{1}{2}(s_i s_j + 1). \quad (14.6)$$

Substituting this expression into Eq. (14.1), we find that

$$Q = \frac{1}{4m} \sum_{ij} B_{ij} (s_i s_j + 1) = \frac{1}{4m} \sum_{ij} B_{ij} s_i s_j, \quad (14.7)$$

where we have used Eq. (14.3) in the second equality. The leading constant factor of $1/4m$ is not important here. It is included by tradition, following the conventional definition of modularity in Eq. (7.54), but it has no effect on the maximization problem.

The quantities B_{ij} in Eq. (14.2) are fixed by the structure of the network. The quantities s_i denote a particular division of the network into two parts. So the modularity maximization problem can be rephrased as follows: given a particular set of values B_{ij} (i.e., a particular network), find the quantities $s_i = \pm 1$ (i.e., the division of the network) that maximize Eq. (14.7).

This problem falls within the generic class of *discrete optimization problems*, problems where we are maximizing the value of a known function of a set of discrete-valued variables. There are a number of general computational methods for finding solutions to such problems, even in cases such as this where the set of possible values of the variables is too large to search through exhaustively. These methods typically only give approximate answers, meaning they can find a solution with a large value of the modularity, but not necessarily the largest overall.

Is this good enough? That depends on the particular network, but in practical situations the answer is often yes. In any network we expect to find, in addition to the overall optimal division, other divisions with modularity close to the optimum but just a little lower. If these suboptimal divisions are similar to the optimal division, in the sense that most nodes are placed in the same groups and only a few are assigned differently, then an approximate modularity maximization scheme will probably work well. Even though it will not, in general, find the overall maximum, if it can reliably return a solution close to the maximum then it will find basically the right division of the network.

However, if there are solutions that have modularity almost as good as the overall maximum but which correspond to very different divisions of the network, in which the assignment of nodes to groups is completely unlike the optimal division, then modularity maximization may fail. It may find a division with a high modularity score, but that division will not be a good guide to the true optimal division.

Quite a lot is known about issues such as these. In particular, it is known that both situations described above can occur. Approximate modularity maxima may or may not be a good guide to the optimal network division, depending on the particular network one is looking at. There is a body of theory, based on ideas from statistical physics concerned with so-called “replica symmetry breaking,” that describes these different regimes of behavior [481]. It is possible, for instance, to create artificial networks to deliberately confound the modularity maximization method [213], networks that by construction have divisions with modularity competitive with the overall maximum that nonetheless correspond to completely different assignments of nodes to groups.

These points may be moot, however, because science is normally interested in the cases where the community structure is clear. If the structure in a network is obscure or difficult to detect, for instance because of competing divisions with similar modularity scores, then in most cases it is not telling us much of interest. Communities that are hard to find simply don’t have much influence on how systems behave. By definition, therefore, we are primarily interested in the easy instances of community detection, the ones where the structure leaps out at us. And for these cases, approximate detection methods usually work well.

14.2.2 A SIMPLE MODULARITY MAXIMIZATION ALGORITHM

Discrete optimization problems have a long history of study in computer science, mathematics, and engineering, and a number of general-purpose techniques have been developed that can be applied to any optimization problem. Among these, several have been applied to modularity maximization, such as *simulated annealing* [227, 329], *genetic algorithms* [299], and *extremal optimization* [151]. In general these methods return results of high quality, but they can be quite slow, making them suitable mainly for applications to smaller networks, up to thousands of nodes, say, or perhaps a little more. For larger networks these methods are not practical.

As an alternative, therefore, researchers have developed a number of methods specific to the modularity maximization problem. We will look at three of these here, starting in this section with a simple node-moving algorithm proposed in Ref. [359]. There are variants of this algorithm that work for divisions of a network into any number of groups, but for the moment let us continue to focus on the two-group case.

The algorithm starts by dividing the network into two equally sized groups at random. Then it considers each node in the network in turn and calculates how much the modularity would change if that node were moved to the other group. In terms of Eq. (14.7), this is equivalent to momentarily flipping the

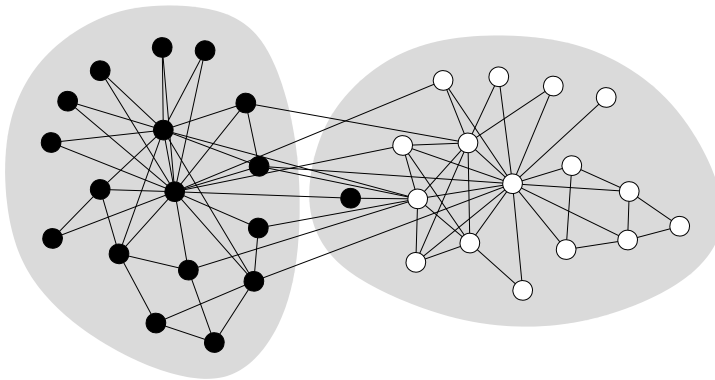


Figure 14.3: Modularity maximization applied to the karate club network. When we apply the node-moving modularity maximization algorithm described in the text to the karate club network, the best division found is the one indicated here by the two shaded regions, which split the network into two groups of 17 nodes each. This division is very nearly the same as the actual split of the network in real life (solid and open circles), following the dispute among the club’s members. Just one node is classified incorrectly.

sign of each of the variables s_i in turn and calculating the effect on Q . The algorithm then chooses from among the nodes the one that most increases, or least decreases, the modularity and it moves this node. Then it repeats the process, but with the important constraint that a node once moved cannot be moved again, at least on this round of the algorithm. When all nodes have been moved exactly once the current round ends.

We then go back over the divisions through which the network passed during the round and select the one with the highest modularity. We use this division as the starting point for another round of the same process, moving nodes one by one as before. We keep repeating this whole procedure, through as many rounds as are needed until the modularity no longer improves, then we stop. The division with highest modularity encountered on the last round is our best estimate of the community structure in the network.

Figure 14.3 shows an example application of this algorithm to the “karate club” network of Zachary [479], which we encountered previously in Chapter 1 (see Fig. 1.2 on page 5). This network represents the pattern of friendships between members of a karate club at a US university, as determined by direct observation of the club’s members over an extended period. The network is interesting because during the period of observation a dispute arose among the members of the club over whether to raise the club’s fees and as a result the club

eventually split into two parts, of 18 and 16 members respectively, the latter departing to form their own club. The colors of the nodes in Fig. 14.3 denote the members of the two factions, while the shaded regions denote the communities identified by the algorithm. As we can see, the communities correspond almost perfectly to the real-life factions. Just one node on the border between the two groups is incorrectly assigned. Thus our algorithm appears to have picked out structure of genuine sociological interest from an analysis of network data alone. It is precisely for results of this kind, that shed light on potentially important structural features of networks, that community detection methods are of interest.

As discussed in Section 8.2, an important feature of any computer algorithm is its computational complexity—how its running time depends on the size of the problem it is solving. The time-consuming part of this algorithm is the calculation of the change in the modularity when a node is moved from one group to the other. Consider the sum in Eq. (14.7) and let us suppose that we are moving node v from one group to the other, so that the variable s_v changes sign while all others remain the same. Only the terms in the sum that contain s_v change their value, so let us write out these terms on their own thus:

$$B_{vv}s_v^2 + s_v \sum_{i(\neq v)} B_{iv}s_i + s_v \sum_{j(\neq v)} B_{vj}s_j. \quad (14.8)$$

The first of these terms $B_{vv}s_v^2$ does not change when we change the sign of s_v so we can ignore it. The other two terms are equal, because $B_{ij} = B_{ji}$, so their sum is just $2s_v \sum_{i(\neq v)} B_{iv}s_i$. Then the change in the modularity, Eq. (14.7), when we flip s_v to $-s_v$ is

$$\Delta Q = \frac{1}{4m} \left[-2s_v \sum_{i(\neq v)} B_{iv}s_i - 2s_v \sum_{i(\neq v)} B_{iv}s_i \right] = -\frac{s_v}{m} \sum_{i(\neq v)} B_{iv}s_i, \quad (14.9)$$

where s_v here denotes the value before the node is moved.

The time-consuming part of calculating ΔQ is the evaluation of the sum

$$\begin{aligned} \sum_{i(\neq v)} B_{iv}s_i &= \sum_{i(\neq v)} A_{iv}s_i - \frac{k_v}{2m} \sum_{i(\neq v)} k_i s_i \\ &= \sum_{i(\neq v)} A_{iv}s_i - \frac{k_v}{2m} \sum_i k_i s_i + \frac{k_v^2 s_v}{2m}, \end{aligned} \quad (14.10)$$

where we have made use of the definition (14.2) of B_{ij} in the first equality.

The first term in this expression is simply the sum of the values of the s_i for all nodes i that are neighbors of the node v being moved. For a network stored

in adjacency list form we can run through these neighbors and calculate the sum in time $O(m/n)$ (see Section 8.3.2), which means it takes a total of $O(m)$ time to perform the calculation for all nodes v . The sum in the second term takes time $O(n)$ to evaluate, since it involves summing over all nodes. However, it does not depend on v , so it takes the same value no matter which node we are moving, meaning it only has to be calculated once for all nodes. The third term in (14.10) can be calculated in $O(1)$ time, or a total of $O(n)$ for all n nodes. Thus the total time taken to calculate ΔQ for all n nodes is $O(m + n)$. A complete round of the algorithm involves repeating this calculation n times, once for each node that actually gets moved, and hence the running time for a round is $O(n(m + n))$.

It is not well understood how the number of rounds required for the whole calculation varies with n or m . In typical applications the number is small, maybe five or ten, at least for networks of up to a few thousand nodes. It is certainly possible that this number could increase for larger networks, but how exactly it does so is not known. Overall, this algorithm is a reasonable one, in terms of speed and quality of results, for applications to smaller networks. Moreover, it is simple to understand and implement. But it is not competitive in terms of speed with the fastest current algorithms, which can have time complexity as good as $O(n \log n)$. We will see an example shortly.

14.2.3 SPECTRAL MODULARITY MAXIMIZATION

The simple algorithm of the previous section works quite well, but is not the fastest method of modularity maximization, nor does it return the very best modularity values. In this section we look at a more advanced “spectral” algorithm that makes use of the properties of the eigenvectors of matrices to rapidly find a good approximation to the maximum modularity division of a network.

Let us again consider the division of a network into just two parts (we will consider the more general case later) and again represent such a division by the quantities

$$s_i = \begin{cases} +1 & \text{if node } i \text{ belongs to group 1,} \\ -1 & \text{if node } i \text{ belongs to group 2.} \end{cases} \quad (14.11)$$

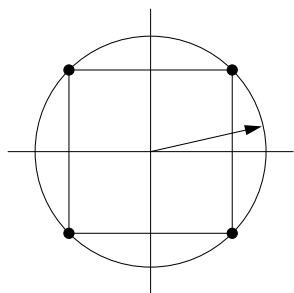
Then the modularity can be written in the form of Eq. (14.7), or alternatively in vector notation as

$$Q = \frac{1}{4m} \mathbf{s}^T \mathbf{B} \mathbf{s}, \quad (14.12)$$

where \mathbf{s} is the n -dimensional vector with elements s_i and \mathbf{B} is the $n \times n$ matrix with elements B_{ij} . The matrix \mathbf{B} is called the *modularity matrix*.

Our goal is to find the division of a given network that maximizes the modularity Q . In other words we wish to find the value of the vector \mathbf{s} that maximizes Eq. (14.12) for a given modularity matrix \mathbf{B} . One reason why this maximization problem is a difficult one is that the s_i can take only a discrete set of values—each s_i is restricted to the values ± 1 . If, on the other hand, the s_i were continuous-valued variables, allowed to take any real values, the problem would be much easier—we could just differentiate to find the maximum.

This suggests the following approximate approach to the maximization problem. Suppose we indeed allow the s_i to take any real values (subject to certain constraints discussed below) and then find the particular values that maximize Q . These values will only give approximately the answer we are looking for, since they probably won't be equal to ± 1 , but they might nonetheless give us a clue about the best division of the network. This idea leads to the so-called *relaxation method*, which is one of the standard methods for the approximate solution of vector optimization problems such as this one. In the present context it works as follows.



The relaxation of the constraint $s_i = \pm 1$ allows the vector \mathbf{s} to point to any position on a hypersphere circumscribing the original hypercube and touching it at the corners.

Note that we have dropped the leading constant $1/4m$ that appears in Eq. (14.7), since, being constant, it has no effect on the position of the maximum.

If we view \mathbf{s} as a vector in an n -dimensional space, then the requirement that $s_i = \pm 1$ implies that the vector must point to one of the corners of an n -dimensional hypercube. Our first step is to relax this constraint and allow the vector to point in any direction in the space. We will, however, keep its length the same. It would not make sense to allow the length to vary—if we did that then we could increase the value of Q , Eq. (14.12), without limit just by making the vector longer and longer, and there would be no maximum to the modularity. So \mathbf{s} will be allowed to take any value, but subject to the constraint that its length stays fixed at the correct value of $|\mathbf{s}| = \sqrt{n}$, or equivalently

$$\sum_i s_i^2 = n. \quad (14.13)$$

Another way of putting this is that \mathbf{s} is now allowed to point to any location on the surface of a hypersphere of radius \sqrt{n} . The hypersphere includes the original allowed values at the corners of the hypercube, but also includes other points in between—see figure.

In this relaxed form the modularity maximization problem is now quite straightforward. We maximize Eq. (14.7) by differentiating, imposing the constraint (14.13) with a Lagrange multiplier β :

$$\frac{\partial}{\partial s_i} \left[\sum_{jk} B_{jk} s_j s_k + \beta \left(n - \sum_j s_j^2 \right) \right] = 0. \quad (14.14)$$

Performing the derivatives we then get

$$\sum_j B_{ij}s_j = \beta s_i, \quad (14.15)$$

or in matrix notation

$$\mathbf{B}\mathbf{s} = \beta\mathbf{s}. \quad (14.16)$$

In other words, the optimal \mathbf{s} is one of the eigenvectors of the modularity matrix and β is the corresponding eigenvalue.

We can work out which eigenvector to use by substituting (14.16) back into Eq. (14.12), which gives

$$Q = \frac{1}{4m}\beta\mathbf{s}^T\mathbf{s} = \frac{n}{4m}\beta, \quad (14.17)$$

where we have used Eq. (14.13), which tells us that $\mathbf{s}^T\mathbf{s} = n$. Since our goal is to make the modularity as large as possible, we want the eigenvalue β to be as large as possible, which means we should choose \mathbf{s} to be the eigenvector corresponding to the largest (most positive) eigenvalue of the modularity matrix.

As we have said, however, the real vector \mathbf{s} is subject to the additional constraint that its elements take the values ± 1 . Typically, this constraint will prevent \mathbf{s} from taking exactly the value given by Eq. (14.16). Let us, however, do the best we can and choose \mathbf{s} to be as close as possible to our ideal value subject to its constraints, which we do by minimizing the angle between \mathbf{s} and the leading eigenvector, which we will denote \mathbf{u} . Equivalently, we can just maximize the inner product $\mathbf{s}^T\mathbf{u} = \sum_i s_i u_i$. The maximum is achieved when $s_i u_i$ is positive for all i , which occurs when s_i has the same sign as u_i for all i :

$$s_i = \begin{cases} +1 & \text{if } u_i > 0, \\ -1 & \text{if } u_i < 0. \end{cases} \quad (14.18)$$

In the unlikely event that a vector element u_i is exactly zero, either value $s_i = \pm 1$ is equally good and we can choose whichever we prefer.

And so we are led to the following very simple algorithm: we calculate the eigenvector of the modularity matrix corresponding to the largest (most positive) eigenvalue and then assign nodes to communities according to the signs of the elements of this vector, positive signs in one group and negative signs in the other.

In practice this method works well. For example, when applied to the karate club network of Fig. 14.3 it works perfectly, classifying every one of the 34 nodes into the correct group.

How fast is it? The most demanding part of the computation is calculating the eigenvector. The leading eigenvector of an $n \times n$ symmetric matrix

can be calculated by methods such as the power method or the Lanczos algorithm [331], which are both based on repeatedly multiplying a vector by the matrix until convergence is reached. The number of multiplications needed depends on the matrix, but for the network problems we are concerned with it typically increases as $\log n$ with network size.⁵ The time needed to perform each multiplication is potentially problematic, given that our matrix \mathbf{B} is not sparse, and indeed usually has all elements non-zero. Normally this would make multiplications rather slow, but by exploiting the particular form of the modularity matrix,⁶ it turns out that it is possible to perform a multiplication in time $O(m + n)$, so that the entire algorithm runs in time $O((m + n) \log n)$, or $O(n \log n)$ on a sparse network with $m \propto n$.

Overall, this spectral method is one of the better methods for modularity maximization, and it also has the benefit of being simple to implement. Assuming one already has a subroutine or library function to calculate the leading eigenvector of a matrix (a standard tool available in most computer languages), the implementation of the algorithm involves only a few lines of programming and is considerably easier, for example, than the node-moving algorithm of Section 14.2.2.

If one is willing to put up with a significantly more complex algorithm, then there are better methods still of maximizing the modularity. The most widely used technique is the so-called Louvain algorithm, which we study in Section 14.2.5. Before we get to that, however, we need to understand how one goes about dividing a network into more than just the two groups we have so far considered.

⁵Specifically, the number of multiplies is $O(\log n)$ for the common case of networks with small diameter. More precisely, if the network is an “expander graph” with a gap between its largest and second-largest eigenvalues that remains constant in the limit of large n , then at most $O(\log n)$ multiplies are needed.

⁶The modularity matrix can be written in vector notation as $\mathbf{B} = \mathbf{A} - \mathbf{k}\mathbf{k}^T/2m$, where \mathbf{k} is the n -element vector whose elements are the degrees k_i of the nodes. Multiplying this matrix into an arbitrary vector \mathbf{v} then gives

$$\mathbf{B}\mathbf{v} = \mathbf{A}\mathbf{v} - \frac{\mathbf{k}^T\mathbf{v}}{2m}\mathbf{k}.$$

The first term on the right-hand side has elements equal to $\sum_j A_{ij}v_j$, which is just the sum of the values v_j on the neighbors of node i . For a network stored in adjacency list format the neighbors can be found and this sum performed in time $O(m/n)$ (see Section 8.3.2), for a total time of $O(m)$ for all n elements. Meanwhile, the second term just requires us to calculate the inner product $\mathbf{k}^T\mathbf{v}$, then multiply it by \mathbf{k} , both of which take time $O(n)$. So the whole multiplication can be done in time $O(m + n)$.

14.2.4 DIVISION INTO MORE THAN TWO GROUPS

The algorithms of the previous two sections perform a limited form of community detection, the division of a network into exactly two communities. But communities are defined to be the natural groupings of nodes in networks and there is no reason to suppose that there will always be just two of them. There might be two, but there might be more, and we would like to be able to find them whatever their number. Moreover we don't, in general, want to have to specify the number of communities; that number should be determined by the structure of the network and not, as above, by the experimenter.

The general problem of finding any number of communities in a network can also be tackled by modularity maximization. Instead of maximizing modularity over divisions of a network into two groups, we just maximize it over divisions into all numbers of groups. Modularity, after all, is supposed to be largest for the best division of the network, no matter how many groups that division possesses.

There are a number of algorithms that adopt this "free maximization" approach, optimizing modularity directly over any number of communities; we discuss some of them at the end of this section and in the following section. First, however, we discuss a simpler approach which is a natural extension of the methods of previous sections, namely repeated bisection of a network. In this approach, we start by dividing the network into two parts and then we further divide those parts in two, and so on. We continue to subdivide as long as doing so increases the modularity of the overall division of the network. When there is no more increase to be gained, we stop.

One must be careful about how one goes about this process, however. We cannot simply treat the groups found in the initial bisection of a network as smaller networks in their own right and apply our bisection algorithm to those smaller networks, because the modularity of the complete network does not break up into independent contributions from the separate groups. Instead, we must consider explicitly the change ΔQ in the modularity of the entire network upon further bisecting a group g of size n_g . Once again using quantities $s_i = \pm 1$ to denote the division of group g , that change can be written as the difference of the modularity before and afterward, thus:

$$\Delta Q = \frac{1}{2m} \sum_{i,j \in g} B_{ij} \frac{1}{2} (s_i s_j + 1) - \frac{1}{2m} \sum_{i,j \in g} B_{ij}, \quad (14.19)$$

all terms involving nodes outside of g canceling. With a little further manipu-

lation we can simplify this expression thus:

$$\begin{aligned} \Delta Q &= \frac{1}{2m} \left[\frac{1}{2} \sum_{i,j \in g} B_{ij} s_i s_j + \frac{1}{2} \sum_{i,j \in g} B_{ij} - \sum_{i,j \in g} B_{ij} \right] = \frac{1}{4m} \left[\sum_{i,j \in g} B_{ij} s_i s_j - \sum_{i,j \in g} B_{ij} \right] \\ &= \frac{1}{4m} \sum_{i,j \in g} \left[B_{ij} - \delta_{ij} \sum_{k \in g} B_{ik} \right] s_i s_j = \frac{1}{4m} \sum_{i,j \in g} B_{ij}^{(g)} s_i s_j, \end{aligned} \quad (14.20)$$

where we have made use of the fact that $s_i^2 = 1$, and we have defined

$$B_{ij}^{(g)} = B_{ij} - \delta_{ij} \sum_{k \in g} B_{ik}. \quad (14.21)$$

Since Eq. (14.20) has the same general form as Eq. (14.7) it can be maximized using the same techniques: the node-moving algorithm of Section 14.2.2, our spectral approach, or any other modularity maximization method will all work for this quantity, just as they did for Eq. (14.7).

The complete algorithm then involves starting with a single division into two groups and repeatedly subdividing using Eq. (14.20), dividing each group again and again, as many times as is necessary. As we have said, since our goal is to maximize the modularity for the entire network, we should go on dividing groups as long as doing so results in an increase in Q . If we are unable to find any division of a group that results in a positive change ΔQ , then we should simply leave that group undivided. When no group can be divided any further, the algorithm is finished.

This repeated bisection method works well in many situations, though it is not perfect. In particular, there is no guarantee that the best division of a network into, say, three parts, can be found by first finding the best division into two parts and then subdividing one of the two. Consider for instance the simple network shown in Fig. 14.4a, which consists of eight nodes in a line. The maximum modularity division of this network into two parts cuts it right down the middle, splitting the network into equally sized groups of four nodes each. The best modularity if the number of groups is unconstrained, however, is that shown in Fig. 14.4b, with three groups of sizes 3, 2, and 3, respectively. A repeated bisection algorithm that finds the optimal two-group division on every step would never find the division in 14.4b because, having first performed the bisection in 14.4a there is no further bisection that will get us to 14.4b.

An alternative approach for dividing networks into more than two communities is to attempt to directly maximize the modularity over divisions into any number of groups. Any of the general-purpose optimization methods discussed at the beginning of Section 14.2.2 can be applied to this multi-group

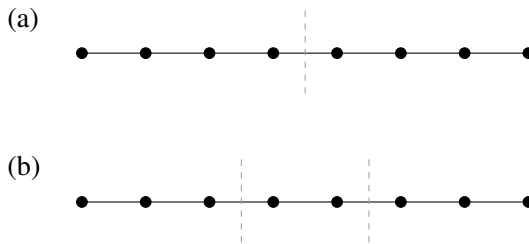


Figure 14.4: Division of a simple network by modularity maximization. (a) The optimal two-group division of this network of eight nodes and seven edges is straight down the middle. (b) The optimal division into an arbitrary number of groups is this division into three.

modularity maximization problem, such as simulated annealing for instance. The node-moving algorithm of Section 14.2.2 also has a natural generalization to the multi-group case: at each step of the algorithm we perform the single move of a node from one group to any other that most increases (or least decreases) the modularity, subject to the constraint that each node is only moved once. The only change from the two-group version of the algorithm is that nodes can be moved to any group. A slight variant that can improve the speed of the calculation is to only consider moves of a node into communities that contain at least one of the node’s neighbors, on the assumption that the other communities are unlikely to be favored anyway.

The spectral algorithm of Section 14.2.3 can also be generalized to the multi-group case [484], although the generalization is not entirely straightforward and involves additional approximations that are not made in the two-group version. By far the most widely used multi-group modularity optimization method, however, is the so-called Louvain algorithm, which we describe next.

14.2.5 THE LOUVAIN ALGORITHM

The Louvain algorithm [66], named after its inventors’ home town in Belgium, is a heuristic algorithm for approximately maximizing modularity over divisions of a network into any number of communities. The algorithm has become a popular choice for many applications because of its speed and because it is included in a number of standard software packages for network analysis, such as Gephi.

The Louvain algorithm is an *agglomerative* algorithm, which works by taking single nodes and joining them into groups, then joining groups with other

groups, and so forth, in an effort to find the configuration with highest modularity. Initially, each node is placed in a separate group on its own. Then one performs a node-moving procedure akin to that of Section 14.2.2, although not identical. In this procedure, one goes through each node in turn and moves that node to another group chosen such that the modularity of the complete system is increased by the largest amount. If no move increases the modularity then the node stays where it is. Also, to make things faster, one only ever considers moving a node into a group that contains at least one of its neighbors. When all nodes have been considered and potentially moved, one repeats the process, and continues to do so until there are no more moves that increase the modularity. This ends the first round of the algorithm.

On the next round, one carries out the same procedure again, but now instead of moving nodes, one moves whole groups. That is, one treats the groups found on the previous round as the units of the algorithm and moves them, in their entirety, from group to group in an effort to increase the modularity, stopping when no further increase is possible.

And so the algorithm proceeds, through as many rounds as are necessary until one reaches a configuration where there are no moves at all that will increase the modularity. This final configuration is then taken as the community division of the network.

The primary advantage of the Louvain algorithm over the others we have discussed is its speed. The exact time complexity of the algorithm is not known, but, like the node-moving algorithm of Section 14.2.2, it appears that a single round typically takes time $O(m + n)$ (because one must consider each node in turn, and each neighbor of that node) and the number of rounds is roughly $O(\log n)$ (because the sizes of the groups roughly double on each round, so at most $\log_2 n$ rounds are possible before the groups reach the size of the whole network). Thus the running time is about $O((m + n) \log n)$, or $O(n \log n)$ on a sparse network with $m \propto n$, which is as good as the spectral algorithm of Section 14.2.3 but doesn't require repeated bisection to find multi-group divisions of a network. In practice, the algorithm is fast enough for very large networks. Its inventors reported one application to a network of over 100 million nodes and a billion edges, with the calculation taking a little over two hours to complete [66].

14.2.6 RESOLUTION LIMIT FOR MODULARITY MAXIMIZATION

The algorithms of the previous sections work well in practice and are widely used. They are, however, not perfect. As pointed out earlier, perfect maximization of the modularity is not practical for any but the smallest of networks, so

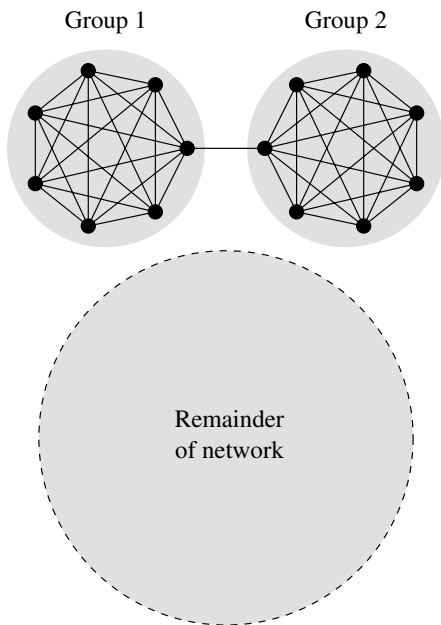


Figure 14.5: Two dense communities within a larger network. The two groups of nodes at the top are indisputably communities by the traditional definition: they are as dense as possible, being complete cliques, while their other connections are very sparse—only a single edge between the two groups and none at all to the rest of the network.

in practice all algorithms are only approximate. Even were it possible to maximize modularity exactly, however, the resulting algorithm would not always give perfect answers, because the modularity maximization method itself is imperfect. Specifically, it suffers from a *resolution limit*, the inability to see communities in a network if they are too small, relative to the size of the network as a whole [184].

Consider the situation depicted in Fig. 14.5, of two communities within a larger network. The two have the same number of nodes, and in this case both are cliques—every possible edge is present within each group. Moreover, neither has any connections to the remainder of the network and there is only a single edge connecting the two groups together, which is the weakest possible connection they could have, other than no connection at all.

Few would dispute that these two groups of nodes are communities in the traditional sense—they could hardly be any denser than they are nor more isolated from the rest of the network and each other. And yet, as we now show, if conditions are right the modularity maximization method will not detect them as separate communities, but will instead erroneously join them together into one.

To demonstrate this, let us calculate the change in the modularity if the two groups are joined. Recall that the modularity is defined to be the fraction

of within-group edges minus the expected fraction of such edges when edge positions are randomized. Thus the change in modularity is equal to the change in within-group edges minus the change in expected edges. When our two groups are joined the number of within-group edges simply increases by 1—the single edge that connects the two becomes a within-group edge. Meanwhile the increase in the expected number of within-group edges is equal to the expected number of edges that fall between the two groups. Let κ_1 and κ_2 be the sums of the degrees of the nodes in each of the two groups, or equivalently the number of “stubs” or ends of edges within each group. In the specific case of our two cliques it is easy to see that $\kappa_1 = \kappa_2 = s^2 - s + 1$, where s is the size of a clique, but let us keep the notation general for the moment and write just κ_1, κ_2 . If edge positions are randomized in the standard fashion in which every stub is equally likely to end up connected to every other (see Section 7.7.1), then the expected number of edges between the two groups will be $\kappa_1\kappa_2/2m$. Taking the difference of the changes in actual and expected numbers of edges between groups, we then find the change ΔQ in modularity upon joining our groups to be

$$\Delta Q = \frac{1}{2m} \left(1 - \frac{\kappa_1\kappa_2}{2m} \right), \quad (14.22)$$

where the initial factor of $1/2m$ is the same overall multiplier that appears in the definition of the modularity, Eq. (14.1).

If this change ΔQ is positive—if the modularity increases upon joining the two groups together—then any algorithm that maximizes modularity correctly will join them. This happens when $1 - \kappa_1\kappa_2/2m > 0$ or equivalently when

$$\kappa_1\kappa_2 < 2m. \quad (14.23)$$

In other words, modularity maximization will fail to distinguish these two groups as separate communities if the product of the sums of their degrees is less than twice the number of edges in the entire network. In a network with 5000 edges, for instance, the method will be unable to distinguish two communities whose degrees each sum to less than 100.

For the specific case of two cliques depicted in Fig. 14.5, we have $\kappa_1 = \kappa_2 = s^2 - s + 1$ where s is the size of a clique, and (14.23) tells us that modularity maximization will fail if $(s^2 - s + 1)^2 < 2m$, or roughly speaking if $s < (2m)^{1/4}$. For instance, in a network with 5000 edges we would not be able to detect clique-like communities of size less than about 10 nodes.

In practice, the resolution limit is not usually a problem for small networks. Community sizes in networks of a few hundred nodes or less rarely approach the limit set by (14.23). For larger networks, however, it can become a problem.

In particular, note that our ability to detect small communities depends on the total number m of edges in the whole network, not just the number in the communities of interest, with the job becoming harder as m increases. Thus, even if the communities themselves do not change, we may lose our ability to detect them just because the network as a whole grows larger.

Consider, for instance, the task of finding communities within a social network. If the social network we are looking at is that of, say, a school with 500 students, then we may well be able to accurately pick out communities within that school using modularity maximization. But if we are looking at the social network of the entire town to which the school belongs, we may find that the very same communities can no longer be detected. Nothing about the communities themselves has changed. They are exactly as they were. But the number of edges m in the network has increased so that the inequality (14.23) becomes satisfied and the modularity maximization method fails.

While modularity maximization is a useful and widely used method, one should bear this issue in mind. Particularly in large networks, you might fail to see small communities even if they are indisputably present in the data.

14.3 METHODS BASED ON INFORMATION THEORY

A completely different approach to community detection is to make use of concepts from information theory, the branch of computer science that deals with measures of information content, for instance in text or numbers. The idea behind these methods is that a good division of the nodes of a network into groups or communities tells us a lot about the structure of the network itself. For instance, if there are more edges within groups than between them then a knowledge of the community structure tells us in which regions of the network most edges will be found. One way to identify good community divisions is thus to search for the divisions that tell us most. Information theory provides tools we can use to do this in a quantitative manner.

Information theory, as it is commonly formulated, deals with linear signals or messages, meaning strings of characters of some kind. In the simplest and most common case there are just two characters, usually denoted 0 and 1, and the messages are binary bit-strings, like 01101011. To apply information theory to networks, one needs to come up with a way of capturing structural features of networks in the form of bit-strings. One way of doing this has been proposed by Rosvall and Bergstrom [416] and is based on the behavior of random walks.

Random walks on networks were discussed previously in Section 6.14.3. Starting at any node in the network, we take a step to one of that node's neighbors chosen uniformly at random, then we repeat the process for as many

We also encountered random walks in Section 4.7 on snowball sampling and in Section 7.1.7 in the context of random-walk betweenness.

steps as we like. (The walk should start at a node with degree greater than zero, so that it has somewhere to go, and ideally at a node within the giant component of the network.) Rosvall and Bergstrom argue that random walks are a natural thing to consider since we are often concerned with flow processes, such as traffic on the Internet, web surfers on the Web, or energy flow in a food web. A random walk is a simple process that captures some of the effects of network structure on flows.

So consider a random walk on an undirected network. The sequence of nodes visited by the walk certainly gives us some information about the structure of the network. At the simplest level it tells us a subset of the edges in the network, since any two nodes that appear consecutively in the sequence must be joined by an edge. However, a random walk also contains more subtle information, information about community structure. Since there are many edges inside communities but few edges between them, a random walk on a network with strong community structure will tend to linger inside communities: when there are few between-group edges along which to escape to another community, the walk will take a long time to find those edges.

In order to quantify the information content of a random walk, Rosvall and Bergstrom turn the walk into a bit-string, a unique string of zeros and ones that exactly describes the walk, as shown in Fig. 14.6. We consider a possible division of the network into communities and we give two labels to each community, an “entry label” and an “exit label,” which take the form of short strings of bits. Every time the random walk enters a new community we record the corresponding entry label. When we leave that community again we record the exit label. Movement of the random walk within a community is recorded in a similar way, by assigning binary labels to the nodes and then recording the labels of the nodes the walk passes through. The complete set of labels, entry and exit labels and node labels, in order, forms the bit-string representation of the random walk.

Rosvall and Bergstrom’s method, which goes by the name of *InfoMap*, focuses on the length of this bit-string. Their central hypothesis is that the best division of the network into communities is the one that corresponds to the shortest bit-string. Clearly the length depends on the lengths of the labels on the communities and nodes, so InfoMap first looks for the set of labels that give the shortest bit-string for a particular community division. That is, keeping the division fixed, one considers all possible ways of assigning the labels to the groups and nodes and chooses the one that gives the shortest overall length of the bit-string, while still representing every possible walk uniquely and unambiguously. Finding the best set of labels is a classic problem in information theory, related to the problem of optimally compressing information such as a

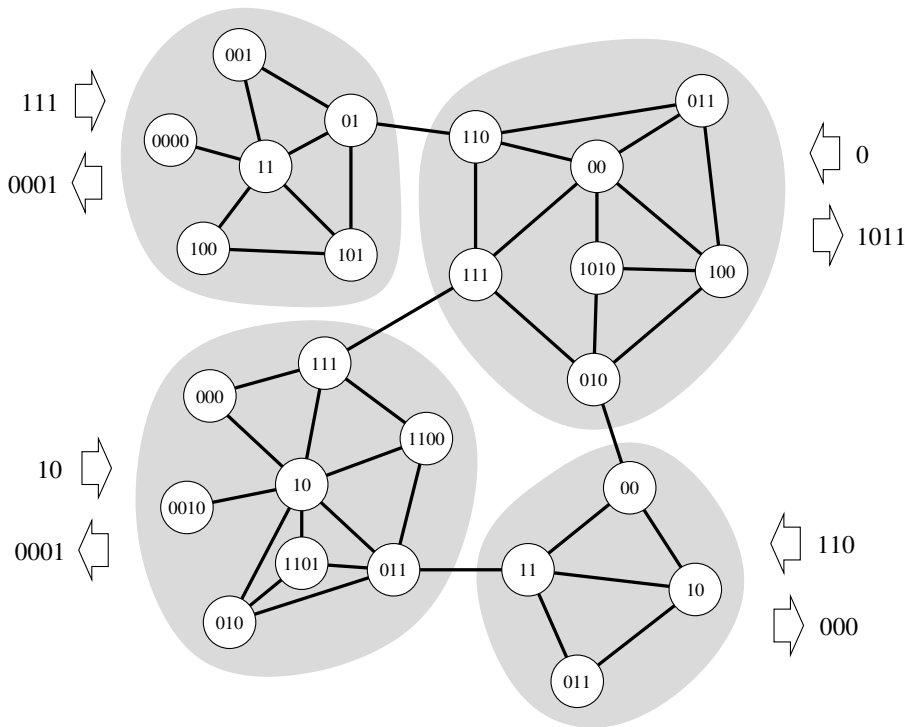


Figure 14.6: Labeling of the nodes and groups in a network with four communities. In the InfoMap method, each node is given a label composed of a short string of zeros and ones. Labels are unique within communities, but nodes in different communities can have the same label. The groups themselves are denoted with entry and exit labels, indicated by the arrows pointing in and out of the groups. A random walk across the network can be uniquely encoded by the sequence of labels of the nodes it visits, along with the entry and exit labels of the groups it enters and leaves. After Rosvall and Bergstrom [416].

file on a computer: the assignment of labels to groups and nodes is a “code” in information theory terms, and we are looking for the optimal code to represent our particular random walk. There are two important points to notice here. First, the labels on the nodes within a community can be the same as the labels in another community. There will be no ambiguity about the random walk if nodes in different communities have the same label, since we always know what group we are in because of the entry and exit labels. The exit labels of different communities can be the same too. The exit label only needs to be distinct from the node labels within the same community, and not from other exit labels (or from node labels in other communities). Second, the node labels

within a community do not all need to have the same length, and in general it makes sense to use shorter labels for nodes that are visited more often. There are fewer short labels than long ones, but if assigned to nodes that are visited frequently their use can still result in significant economy in the length of the final bit-string.

But the length of the bit-string also depends on the specific community division we assume, and in particular it will be short for “good” divisions and long for “bad” ones. There are two reasons for this. First, for a good division, one in which there are few edges between groups, the random walk will, as we have said, move between groups only rarely, and hence we will not have to record the entry and exit labels of groups very often. At the same time, if the groups are reasonably small then the node labels within groups can be kept short. There are 2^b distinct labels with b bits, so in a group of N nodes we can label them all distinctly with labels of $b \simeq \log_2 N$ bits. Again, however, we can economize by using shorter labels for frequently visited nodes and hence may be able to do better than $\log_2 N$ on average.

In some sense the group and node label processes are at odds with one another: choosing larger groups will make transitions between groups rarer, but smaller groups allow us to use shorter group labels. If we want to reduce the total length of the bit-string, therefore, we will have to compromise between having a few large groups and many small ones. The assumption of the InfoMap method is that the best compromise is given by the coding of the random walk that results in the shortest bit-string.

One can thus imagine performing the following steps. First, we generate a random walk long enough to visit all parts of the network. Then, for each possible division of the network into groups, we find the set of group and node labels that gives the shortest bit-string representation of that random walk. Then we look through all divisions to find the one that gives the shortest bit-string overall, and this is deemed to be the best community structure for the network.

This method would work in theory, but it would be very laborious, and in practice things can be done much faster by making use of some results from information theory. The core of the calculation is finding the length of the shortest bit-string representation of the random walk for a particular community division. A fundamental result, known as Shannon’s source coding theorem, tells us that for the shortest possible bit-string the average number L of bits per step of the random walk is equal to the *entropy* of the random walk,

which in this case is given by the *map equation*:

$$L = qH(Q) + \sum_g p_g H(P_g). \quad (14.24)$$

Here q is the fraction of the time that the random walk spends hopping between groups, and p_g is the fraction of the time it spends within group g and exiting group g . The quantities $H(Q)$ and $H(P_g)$ are information-theoretic entropies. The entropy of a sequence Q of objects is given by

$$H(Q) = - \sum_i Q_i \log_2 Q_i, \quad (14.25)$$

where Q_i is the fraction of times that object i appears in the sequence. In Eq. (14.24) $H(Q)$ is the entropy of the sequence of groups that the random walk passes through (strictly the sequence of entry labels) and $H(P_g)$ is the entropy of the nodes in group g that the walk passes through (strictly the sequence of node labels and the exit label—the exit label is considered part of the group for this purpose).

The map equation allows us to calculate L without actually assigning any labels to groups or nodes: the group and node labels are a useful thought experiment to motivate the method, but in the end they do not enter the calculation. Moreover, it turns out we don't actually need to perform a random walk either. One can calculate the probabilities p_g and q and the entropies $H(P_g)$ and $H(Q)$ simply by knowing the local structure of the network and the fraction of the time the random walk spends at each node, which is simply proportional to the degree of the node. Putting everything together, one can then calculate L rapidly and minimize it over community divisions.

As in the case of modularity maximization, there are in practice too many possible divisions to search through them exhaustively for the smallest value of L in any but the smallest of networks, so one must make use of heuristic optimization strategies. In their work, Rosvall and Bergstrom used an algorithm similar to the Louvain method of Section 14.2.5, minimizing L by moving individual nodes from group to group, then repeating the exercise at the level of entire groups, and so on until no further improvement is possible.

In the end, therefore, the InfoMap method is in some ways rather similar to modularity maximization: it defines a quality function, in this case the entropy L , which characterizes how good a particular community division is, then optimizes it over possible divisions to find the best one. The motivation behind the quality function is quite different from the motivation behind the modularity, but both methods ultimately reduce to an optimization problem.

The InfoMap method appears to work very well, returning high-quality results in standardized tests (see Section 14.6). It is also fast. Although there are no formal results for its time complexity, it appears to have a running time similar to that of the ordinary Louvain algorithm, i.e., about $O(n \log n)$ on a sparse network. It is a somewhat complex algorithm to implement, but good implementations are available to save you the effort, and the algorithm has become a favorite due to its excellent results.

14.4 METHODS BASED ON STATISTICAL INFERENCE

Some of the most powerful and flexible methods for community detection are those based on *statistical inference*. These methods work by fitting a network model—typically some kind of random graph—to observed network data. The parameters of the fit can tell us about features of the network, including community structure, in much the same way that the fit of a straight line through a set of data points can tell us about their slope.

See Chapters 11 and 12 for a discussion of random graph models.

We encountered statistical inference methods, and specifically the method of maximum likelihood, previously in Chapter 9 on error estimation. In Section 9.3.2 we studied the simple (non-network) example of a set of data points drawn from a normal or Gaussian distribution and showed that given only the points themselves, plus the knowledge that they are normally distributed, we can estimate the parameters of the Gaussian—the mean μ and standard deviation σ —by writing down an expression for the probability, or likelihood, of the data and then maximizing it.

In the language of statistics, the Gaussian distribution in this example is a “model.” It might not fit our usual notions of what a model should be, but to a statistician a model is any process that can generate data. Most models have one or more free parameters that characterize them (such as the parameters μ and σ in this case) and by maximizing the likelihood we can find the values of the parameters that give the best match between the model and a given set of data.

The same exact approach can also be applied to networks. Given a network model, meaning any process that can generate a network, we can fit that model to data—i.e., to a particular network structure—by finding the values of the model parameters that give the highest likelihood. In effect, we are saying, “If this network was generated by this model, what is our best guess at the values of the model parameters that were used?” As we will see, such fits of models to data can often shed a lot of light on network structure.

As an example, consider the Poisson random graph model of Chapter 11. Recall that, other than the size of the network n , this model has just one

parameter, the probability p that any two distinct nodes are connected by an edge. Every edge is independent and has the same probability, so the total probability—the likelihood—that a particular network, defined by its adjacency matrix \mathbf{A} , is generated by the random graph model with a particular value of p , is

$$P(\mathbf{A}|p) = p^m (1-p)^{\binom{n}{2}-m}, \quad (14.26)$$

where m is the number of edges in the network as usual. In other words there are $\binom{n}{2}$ pairs of nodes in total, and we get a factor of p for each of the m pairs that are connected by an edge, and a factor of $1-p$ for each of the $\binom{n}{2} - m$ pairs that are not connected.

Now suppose that we do not know the value of p . All we are given is the data, i.e., the network itself. We can make an estimate of p by employing Bayes' rule thus:

$$P(p|\mathbf{A}) = P(\mathbf{A}|p) \frac{P(p)}{P(\mathbf{A})}, \quad (14.27)$$

where $P(p)$ and $P(\mathbf{A})$ are the prior probabilities on p and \mathbf{A} respectively. The most probable value of p is now, by definition, given by maximizing this expression with respect to p while holding \mathbf{A} constant at its observed value. But if \mathbf{A} is constant then so is $P(\mathbf{A})$, meaning that the denominator in (14.27) has no effect on the position of the maximum. Moreover, we typically also assume that $P(p)$ is constant as well, i.e., that all values of p from zero to one are equally likely a priori. With this assumption, maximizing $P(p|\mathbf{A})$ is equivalent to maximizing the likelihood $P(\mathbf{A}|p)$, so far as determining the value of p goes. They both have their maximum in the same place.

The maximization can be performed in this case by simple differentiation. Differentiating Eq. (14.26) and setting the result to zero gives

$$mp^{m-1}(1-p)^{\binom{n}{2}-m} - \left[\binom{n}{2} - m\right]p^m(1-p)^{\binom{n}{2}-m-1} = 0, \quad (14.28)$$

which can be rearranged to read

$$p = \frac{m}{\binom{n}{2}}. \quad (14.29)$$

As discussed in Section 9.3.2, we often work not with the likelihood itself but with its logarithm, which in this case has the value

$$\log P(\mathbf{A}|p) = m \log p + \left[\binom{n}{2} - m\right] \log(1-p). \quad (14.30)$$

Since the logarithm is a monotone increasing function of its argument, its maximum is in the same place as the maximum of the likelihood itself, but

differentiating the logarithm is often algebraically simpler. We will use this trick in some of our later calculations. We leave it as an exercise to verify that differentiating (14.30) does indeed give Eq. (14.29) again.

In a sense Eq. (14.29) is trivial. It says that the best estimate of the edge probability p is the obvious one—the fraction of the $\binom{n}{2}$ pairs of nodes that are actually connected by edges. Still, it's good to see that the maximum likelihood method does give sensible results when applied to networks. Let us now look at a less trivial application.

14.4.1 COMMUNITY DETECTION USING STATISTICAL INFERENCE

We can use the maximum likelihood method to perform community detection by fitting network data to a model that contains community structure. The model we use is the degree-corrected stochastic block model introduced in Section 12.11.6. Recall that in this model we divide the nodes of a network into some number q of groups or communities, labeled by integers $1 \dots q$, and then place undirected edges between node pairs with probability $\omega_{g_i g_j} c_i c_j / 2m$, where g_i and g_j are the groups to which nodes i and j belong and c_i is the desired average degree of node i . The $q \times q$ matrix of parameters ω_{rs} controls the community structure. For instance, if the diagonal entries ω_{rr} are larger than the off-diagonal ones, the network will have traditional assortative community structure in which connections are more likely within groups than between them. The model can, however, capture other kinds of structure as well, such as disassortative structure where the diagonal entries are smaller than the off-diagonal ones.

We might imagine we could use the non-degree-corrected version of the stochastic block model also, but it turns out this gives poor results. See Ref. [257].

As described in Section 12.11.6, if the c_i are to be equal to the degrees of the nodes on average, then the parameters ω_{rs} need to satisfy an additional set of q constraints

$$\sum_j \omega_{r g_j} c_j = 2m, \tag{14.31}$$

one for each value of r (see Eq. (12.143)). For our purposes here, it will be convenient to rewrite the left-hand side of this equation thus:

$$\sum_j \omega_{r g_j} c_j = \sum_{j_s} \omega_{rs} \delta_{g_j s} c_j = \sum_s \omega_{rs} \kappa_s, \tag{14.32}$$

where δ_{ij} is the Kronecker delta and we have defined the quantity

$$\kappa_s = \sum_j \delta_{g_j s} c_j, \tag{14.33}$$

which is just the sum of the average degrees c_j of all nodes in group s . Combining Eqs. (14.31) and (14.32), we now have

$$\sum_s \omega_{rs} \kappa_s = 2m. \quad (14.34)$$

This completely defines the degree-corrected stochastic block model. In fact, however, when used for community detection the model is usually studied in a slightly different version in which, rather than just placing a single edge between any pair of nodes, we place a Poisson-distributed number of edges with mean equal to $\omega_{g_i g_j} c_i c_j / 2m$, or a half of this value when $i = j$. This allows there to be multiedges in the network, as there are, for instance, in the configuration model also—see Section 12.1. As with the configuration model the presence of multiedges is not realistic for many networks, but typically the value $\omega_{g_i g_j} c_i c_j / 2m$ is a small number because m is large, so the probability of having two or more edges between a pair of nodes is very small—almost all pairs of nodes will have either one edge or none, so multiedges can be neglected. Self-edges are also allowed in the network, but they too will be few in number. Thus using this version of the model makes little difference to the networks generated, but it turns out to make the calculations significantly easier.

The degree-corrected stochastic block model has three sets of parameters. The first is the $q \times q$ matrix with elements ω_{rs} , which we will denote $\mathbf{\Omega}$. Note that this matrix is symmetric, since the probability ω_{rs} of an edge between group r and group s is necessarily equal to that for an edge between s and r . The second set of parameters is the quantities c_i , equal to the average degrees of the nodes, which we can think of as a vector \mathbf{c} with n elements, one for each node.

There is also a third, hidden set of parameters that you might not notice at first, namely the groups g_i to which the nodes belong, which we will consider as the elements of a vector \mathbf{g} . To specify the model completely one must specify these group assignments, so properly they should be regarded as parameters too. It is by calculating the best-fit values of these parameters that we will perform community detection using the maximum likelihood method.

Having defined our model, we can now write down the likelihood that a network with adjacency matrix \mathbf{A} is generated by the degree-corrected stochastic

block model:

$$\begin{aligned}
 P(\mathbf{A}|\mathbf{\Omega}, \mathbf{c}, \mathbf{g}) &= \prod_{i<j} \frac{(\omega_{g_i g_j} c_i c_j / 2m)^{A_{ij}}}{A_{ij}!} e^{-\omega_{g_i g_j} c_i c_j / 2m} \\
 &\quad \times \prod_i \frac{(\omega_{g_i g_i} c_i^2 / 4m)^{A_{ii}/2}}{(\frac{1}{2}A_{ii})!} e^{-\omega_{g_i g_i} c_i^2 / 4m}. \quad (14.35)
 \end{aligned}$$

There are several points to note about this expression. First of all, it is, at heart, just a product of Poisson distributions, one for each pair of nodes i, j , that represent the probability of observing specific values A_{ij} of the adjacency matrix elements. Notice how we have broken the expression into separate products for self-edges and non-self-edges, because the two have slightly different forms. The first product represents the non-self-edges. Taking the product $\prod_{i<j}$ over node pairs with $i < j$ ensures both that this term contains no self-edges $i = j$ and that each distinct pair of nodes is counted only once, not twice.

The second product accounts for the self-edges and differs from the first in a couple of ways. First, note the appearance of $\frac{1}{2}A_{ii}$, where in the first product we had A_{ij} . Recall that, conventionally, a self-edge in a network is represented by setting $A_{ii} = 2$ (see Section 6.2), so $\frac{1}{2}A_{ii}$ correctly gives the number of self-edges at node i . Second, note that the mean of the Poisson distribution in the self-edge term is $\omega_{g_i g_i} c_i^2 / 4m$. The 4 in the denominator is because, as defined above, the mean number of self-edges at a node in this model is a half of the corresponding number for non-self-edges. Alternatively, $\omega_{g_i g_j} c_i c_j / 2m$ is the mean value of every adjacency matrix element, both diagonal and off-diagonal, but the diagonal elements A_{ii} are equal to twice the number of self-edges, and hence the mean number of self-edges is $\omega_{g_i g_i} c_i^2 / 4m$.

As mentioned previously, it is usually easier to work with the logarithm of the likelihood than with the likelihood itself. Taking the log of Eq. (14.35), we get

$$\begin{aligned}
 \log P(\mathbf{A}|\mathbf{\Omega}, \mathbf{c}, \mathbf{g}) &= \sum_{i<j} \left[A_{ij} \log \frac{\omega_{g_i g_j} c_i c_j}{2m} - \log A_{ij}! - \frac{\omega_{g_i g_j} c_i c_j}{2m} \right] \\
 &\quad + \sum_i \left[\frac{1}{2}A_{ii} \log \frac{\omega_{g_i g_i} c_i^2}{4m} - \log(\frac{1}{2}A_{ii})! - \frac{\omega_{g_i g_i} c_i^2}{4m} \right]. \quad (14.36)
 \end{aligned}$$

Since our goal is to find the maximum of this expression with respect to the model parameters we can ignore constant terms that do not depend on the parameters, such as the term in $\log A_{ij}!$. Gathering together the remaining

terms, Eq. (14.36) can then be simplified to

$$\log P(\mathbf{A}|\mathbf{\Omega}, \mathbf{c}, \mathbf{g}) = \frac{1}{2} \sum_{ij} \left[A_{ij} \log \frac{\omega_{g_i g_j} c_i c_j}{2m} - \frac{\omega_{g_i g_j} c_i c_j}{2m} \right] + \text{constants}. \quad (14.37)$$

Note how the sum now includes terms for both $i < j$ and $i > j$, but we have compensated for this double counting by including an extra factor of $\frac{1}{2}$ in front of the whole expression.

This formula can be simplified a little further. The first term of the sum can be rewritten thus:

$$\begin{aligned} \sum_{ij} A_{ij} \log \frac{\omega_{g_i g_j} c_i c_j}{2m} &= \sum_{ij} A_{ij} \log \omega_{g_i g_j} + \sum_{ij} A_{ij} \log c_i \\ &\quad + \sum_{ij} A_{ij} \log c_j - \sum_{ij} A_{ij} \log 2m \\ &= \sum_{ij} A_{ij} \log \omega_{g_i g_j} + \sum_i k_i \log c_i + \sum_j k_j \log c_j - 2m \log 2m, \end{aligned} \quad (14.38)$$

where we have made use of the fact that $\sum_j A_{ij} = k_i$ (Eq. (6.12)) and $\sum_{ij} A_{ij} = 2m$ (Eq. (6.13)). Now we notice several things. First, the final term $-2m \log 2m$ is a constant, independent of any of our parameters, and hence can be neglected. Furthermore, the two middle terms $\sum_i k_i \log c_i$ and $\sum_j k_j \log c_j$ are the same—the only name of the summation variable has changed. And the first term we can rewrite, using a trick similar to that in Eq. (14.32), as

$$\sum_{ij} A_{ij} \log \omega_{g_i g_j} = \sum_{ijrs} \delta_{g_i r} \delta_{g_j s} A_{ij} \log \omega_{rs} = \sum_{rs} m_{rs} \log \omega_{rs}, \quad (14.39)$$

where we have defined the quantity

$$m_{rs} = \sum_{ij} \delta_{g_i r} \delta_{g_j s} A_{ij}, \quad (14.40)$$

which is equal to the number of edges that run between groups r and s in our network, or twice that number when $r = s$, since each pair of nodes i, j gets counted twice in that case.⁷

⁷One can think of m_{rs} as a group-level equivalent of the adjacency matrix. Its off-diagonal elements give the number of edges between groups and its diagonal elements give twice the number within groups, analogous to the diagonal elements of the adjacency matrix.

Putting everything together, we now have

$$\log P(\mathbf{A}|\mathbf{\Omega}, \mathbf{c}, \mathbf{g}) = \sum_i k_i \log c_i + \frac{1}{2} \sum_{rs} m_{rs} \log \omega_{rs} - \frac{1}{2} \sum_{ij} \frac{\omega_{g_i g_j} c_i c_j}{2m} + \text{constants.} \quad (14.41)$$

This is the log-likelihood for the degree-corrected stochastic block model. Armed with this expression, we can now calculate the best-fit values of the model parameters by maximizing. In particular, since the group memberships g_i are among our parameters, we can calculate the maximum likelihood assignment of the nodes to the q groups. This is how we perform community detection: we find the assignments g_i of nodes to groups that are most likely given the observed network \mathbf{A} . In this sense, the inference method is a particularly clear and rigorous method of community detection. With the other methods we have looked at one could argue about, for instance, the particular definition of the modularity or whether the information theoretic assumptions of the InfoMap method are justified. But the inference method simply tells us which is the most likely division of the network into groups, given the observed network structure.

The weakness of the method, on the other hand, is that we have assumed that the network was generated using the degree-corrected stochastic block model. Most real-world networks, presumably, are generated by processes different from—and probably more complicated than—the block model, so this assumption is somewhat dubious. Nonetheless, it turns out that the maximum likelihood method applied to the block model gives excellent results in most cases. It is one of the central mysteries not just of this field, but of the whole of modern statistics, that models like this one, which obviously fail to capture many of the intricate details of the real world, nonetheless give good results when applied to real-world data.

In any case, let us press ahead with our calculation. Maximizing the log-likelihood with respect to c_i and ω_{rs} is straightforward—we can just differentiate as usual. Maximizing with respect to the group memberships g_i is harder because they are discrete variables, so differentiating will not work. But let us do the easy steps first.

Differentiating Eq. (14.41) with respect to c_i gives

$$\frac{\partial \log P}{\partial c_i} = \frac{k_i}{c_i} - \sum_j \frac{\omega_{g_i g_j} c_j}{2m} = \frac{k_i}{c_i} - 1, \quad (14.42)$$

where we have made use of the fact that $\omega_{rs} = \omega_{sr}$ in the first equality and Eq. (14.31) in the second. Setting this expression equal to zero and solving for c_i

we find that

$$c_i = k_i. \quad (14.43)$$

In other words, the best choice for the expected degree parameters c_i is just to set them equal to the observed degrees of the nodes in the network.

In differentiating (14.41) with respect to ω_{rs} it will be convenient to rewrite the final term slightly:

$$\sum_{ij} \frac{\omega_{g_i g_j} c_i c_j}{2m} = \sum_{ijrs} \delta_{g_i r} \delta_{g_j s} \frac{\omega_{rs} c_i c_j}{2m} = \sum_{rs} \frac{\omega_{rs} \kappa_r \kappa_s}{2m}, \quad (14.44)$$

where κ_r and κ_s are the sums of the c_i within groups r and s , as in Eq. (14.33). Then we have

$$\log P(\mathbf{A}|\mathbf{\Omega}, \mathbf{c}, \mathbf{g}) = \sum_i k_i \log c_i + \frac{1}{2} \sum_{rs} \left(m_{rs} \log \omega_{rs} - \frac{\omega_{rs} \kappa_r \kappa_s}{2m} \right) + \text{constants}. \quad (14.45)$$

Differentiating this expression with respect to ω_{rs} and setting the result to zero now gives⁸

$$\frac{m_{rs}}{\omega_{rs}} - \frac{\kappa_r \kappa_s}{2m} = 0, \quad (14.46)$$

or

$$\omega_{rs} = 2m \frac{m_{rs}}{\kappa_r \kappa_s}. \quad (14.47)$$

We should bear in mind that the ω_{rs} are required to satisfy the q constraints of Eq. (14.34). That Eq. (14.47) does in fact do this we can confirm by substituting into (14.34) to get

$$\sum_s \omega_{rs} \kappa_s = \sum_s 2m \frac{m_{rs}}{\kappa_r \kappa_s} \kappa_s = \frac{2m}{\kappa_r} \sum_s m_{rs}. \quad (14.48)$$

Using the definition of m_{rs} in Eq. (14.40), we have

$$\sum_s m_{rs} = \sum_{ijs} \delta_{g_i r} \delta_{g_j s} A_{ij} = \sum_{ij} \delta_{g_i r} A_{ij} = \sum_i \delta_{g_i r} k_i = \sum_i \delta_{g_i r} c_i = \kappa_r, \quad (14.49)$$

where we have used Eq. (14.43). Substituting this result back into (14.48) we then see that (14.34) is satisfied.⁹ Thus Eq. (14.47) correctly maximizes the likelihood while respecting the constraints on ω_{rs} .

⁸In performing the derivative we should bear in mind that $\omega_{rs} = \omega_{sr}$, so that a derivative with respect to one is also a derivative with respect to the other. All this does, however, is introduce an extra factor of two in both the first and second terms of Eq. (14.46), factors that immediately cancel out again, leaving the result as we see it.

⁹It is somewhat fortuitous that the constraints are satisfied in this way. Normally, one would expect to have to enforce the constraints explicitly using Lagrange multipliers. One can in fact do that here, and if one does so it leads to the same result.

Having determined the best-fit values of the parameters c_i and ω_{rs} , we can now substitute these back into the expression for the log-likelihood, Eq. (14.41), to get the so-called *profile likelihood*:

Technically, this is a profile log-likelihood, but it's usually just called the profile likelihood.

$$\mathcal{L} = \frac{1}{2} \sum_{rs} m_{rs} \log \frac{m_{rs}}{\kappa_r \kappa_s} + \text{constants}, \quad (14.50)$$

where several terms have been neglected because they are constant. (For instance, any term that depends only on k_i is constant.)

This expression is the fundamental equation for community detection using maximum likelihood. It tells us the value of the log-likelihood after it has been maximized with respect to the continuous parameters \mathbf{c} and $\mathbf{\Omega}$. All that remains is for us to maximize with respect to the group memberships \mathbf{g} , which enter the expression through the values of m_{rs} , κ_r , and κ_s . For any particular assignment \mathbf{g} of nodes to groups, we can calculate m_{rs} , κ_r , and κ_s from Eqs. (14.33) and (14.40), substitute them into Eq. (14.50) to get a value for \mathcal{L} , then maximize this value over all possible assignments.

Thus, the problem of community detection is, once again, reduced to the maximization of a function over possible divisions of the network. As with the previous methods we have considered, exhaustive maximization over all possible divisions is not practical except for very small networks. There are simply too many divisions to search through them all, so one must employ approximate heuristics. For instance, Newman and Karrer [257] employed a node-moving algorithm that is the equivalent of the modularity maximization algorithm of Section 14.2.2—one repeatedly moves a node from one group to another, choosing at every step the move that most increases, or least decreases, the profile likelihood (14.50). In principle, one could also create an equivalent of the Louvain algorithm of Section 14.2.5, and any generic optimization scheme could be applied to the problem, such as simulated annealing or a genetic algorithm.

This method is found to give excellent results in practice, and moreover is highly regarded because of its principled conceptual foundation, which arguably lends the results greater credence than those derived from more ad hoc methods such as modularity maximization. In terms of running speed it is comparable to modularity maximization or the information theoretic method of Section 14.3, since the algorithms are in practice very similar, differing only in the particular function they are optimizing. Because of its solid mathematical foundations the maximum likelihood method is also amenable to rigorous analysis in a way that some other methods are not. Bickel and Chen [62], for instance, proved that, under suitable conditions, the method is *asymptotically consistent*, meaning that it will correctly identify known community structure

in the limit of large network size.

Conversely, Decelle *et al.* [138, 139], working with a different version of the method, proved that there are other conditions under which it will fail to identify known community structure, if that structure is too weak. They then used this result as a starting point to prove a much more general result, that under the same circumstances *no algorithm* of any kind will be able to detect the communities. This latter result is a particularly powerful one in that it demonstrates that some community detection problems are simply impossible. There are some situations where structure exists in a network but no algorithm will find it.¹⁰

There are some disadvantages to the maximum likelihood method. The main one is that it requires us to specify the number q of groups in the network from the outset, where other methods such as modularity maximization or InfoMap treat q as a free variable and find its optimal value as an integral part of the community detection process. To apply the maximum likelihood approach in the general case where the value of q is not known one needs to use other techniques to estimate q first. Such techniques do exist, typically based on Bayesian model selection [119, 408, 475] or minimum description length methods [385], although for the moment at least they are rather slow and hence applicable only to relatively small networks.

14.5 OTHER ALGORITHMS FOR COMMUNITY DETECTION

We have seen three approaches to community detection in this chapter, based on modularity maximization, on information theoretic methods, and on statistical inference. These methods are some of the most accurate and widely used but they are by no means the only methods that have been proposed. This is a highly active area of research and there are by now an impressive number of different approaches, methods, and algorithms for the community detection problem. In this section we describe briefly some of the best known of these other approaches.

¹⁰On the other hand, if no algorithm can detect the presence of community structure in a network, then no real-world process, physical, social, technological, or otherwise, can have an outcome that depends on it either. If it did, then one could simulate that process to create a computer algorithm that detects the presence of the community structure, which is impossible. Arguably, therefore, we don't care about community structure in regimes where it is undetectable, since it can never affect any process we care about.

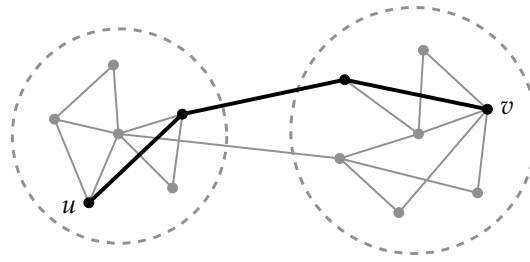


Figure 14.7: Identification of between-group edges. This small network is divided into two groups of nodes (marked by the dashed circles), with only two edges connecting the groups. Any path joining nodes in different groups (such as nodes u and v) must necessarily pass along one of these two edges. Thus if we consider paths between all pairs of nodes we expect the between-group edges to carry more paths than most. By counting the number of paths that pass along each edge we can in this way identify the between-group edges.

14.5.1 BETWEENNESS-BASED METHODS

All of the methods we have seen so far for finding communities involve the optimization of various measures of community structure, such as the modularity of Section 14.2 or the likelihood of Section 14.4. A completely different approach to the problem is to look for the edges in the network that lie between communities. If we can find and remove these edges, we will be left with just the isolated communities.

There is more than one way to quantify what we mean by an edge that lies between communities, but one common approach is to use betweenness centrality. As described in Section 7.1.7, the betweenness centrality of a node in a network is the number of shortest paths in the network that pass through that node. Similarly, we can define an *edge betweenness* that counts the number of shortest paths that run along edges, and edges that lie between communities can be expected to have high values of this edge betweenness—see Fig. 14.7.

The calculation of edge betweenness is closely analogous to the node case: we consider the shortest path or paths between every pair of nodes in the network (except nodes in different components, for which no such path exists), and count how many such paths go along each edge. Edge betweenness can be calculated for all edges in time $O(n(m + n))$ using a slightly modified version of the algorithm described in Section 8.5.6 [366].

Our algorithm for detecting communities is then as follows. We calculate the betweenness scores of all edges in our network, then find the edge with

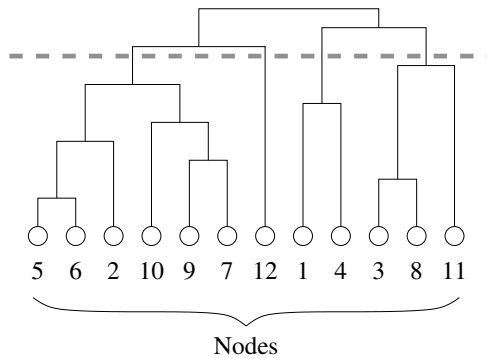


Figure 14.8: A dendrogram. The results of the edge betweenness algorithm can be represented using a tree or “dendrogram,” in which the nodes are depicted (conventionally) at the bottom and the “root” at the top of the tree represents the whole network. The fragmentation of the network as edges are removed one by one is represented by the successive branching of the tree as we move down the figure and the identities of the nodes in a connected subset at any point in the procedure can be found by following the tree down to the bottom. Each intermediate division of the network through which the algorithm passes corresponds to a horizontal cut through the dendrogram. For instance, the cut denoted by the dotted line in this dendrogram splits the network into four groups of 6, 1, 2, and 3 nodes respectively.

the highest score and remove it. In removing this edge we will change the betweenness scores of some other edges, because any shortest paths that previously traversed the removed edge will now have to be rerouted another way. So we recalculate the betweenness scores following the removal, then we search again for the edge with the highest score and remove it, and so forth.

As we remove one edge after another, an initially connected network will eventually split into two pieces, then into three, and so on. The progress of the algorithm can be represented using a tree or *dendrogram* as shown in Fig. 14.8. At the bottom of the figure we have the “leaves” of the tree, which represent the nodes of the network. As we move up the tree the leaves join together first into pairs and then into larger groups, until at the top of the tree all are joined together to form a single whole. Our algorithm in fact generates the dendrogram in the opposite direction, from the top down, starting with a single connected network and splitting it repeatedly until we get to the level of single nodes. Intermediate configurations of the network during the run of the algorithm correspond to horizontal cuts through the dendrogram, as indicated by the dotted line in the figure. Each branch of the tree that intersects this dotted line represents one group of nodes, whose membership we can determine by

following the branch down to its leaves at the bottom of the figure. Thus the dendrogram captures in a single diagram the configuration of groups in the network at every stage from start to finish of the calculation.

This algorithm is thus somewhat different from previous ones, in that it doesn't give a single decomposition of a network into communities, but a set of different possibilities, ranging from coarse divisions into just a few large communities (at the top of the dendrogram) to fine divisions into many small communities (at the bottom). It is up to the user to decide which of the many divisions represented is most useful for their purposes. One could, in principle, use a measure such as modularity to quantify the goodness of the different divisions and select the best one, but this somewhat misses the point. If high modularity is what you care about, then you are better off simply maximizing modularity directly. It is more appropriate to think of the betweenness algorithm as producing a different kind of output, one that has its own advantages and disadvantages but that can undoubtedly tell us interesting things about network structure.

The betweenness algorithm is, unfortunately, quite slow. As we have said, the calculation of betweenness for all edges takes time of order $O(n(m+n))$ and we must perform this calculation repeatedly, once for each edge removed, so the entire algorithm takes time $O(mn(m+n))$ to remove all m edges in a network, or $O(n^3)$ on a sparse network with $m \propto n$. This makes the algorithm one of the slower ones considered in this chapter. It does give quite good results in practice [204,366], but it has mostly been superseded by the faster methods of previous sections.

Nonetheless, the ability of the algorithm to return an entire dendrogram, rather than just a single division of a network, could be useful in some cases. The divisions represented in the dendrogram form a hierarchical decomposition of the network in which the communities at one level are completely contained within the larger communities at all higher levels. There has been some interest in hierarchical structure in networks and hierarchical decompositions that might capture it. We look at another algorithm for hierarchical decomposition in Section 14.5.2.

A variant of the betweenness algorithm has been proposed by Radicchi *et al.* [395] based on the same idea of identifying the edges between communities and removing them, but using a different measure to perform the identification. Radicchi *et al.* observe that the edges that fall between otherwise poorly connected communities are unlikely to belong to short loops of edges, since doing so would require that there be two nearby edges joining the same groups—see Fig. 14.9. Thus one way to identify the edges between communities is to look for edges that belong to an unusually small number of

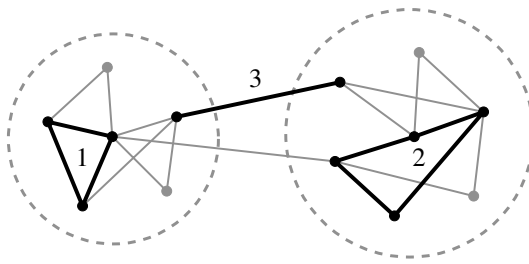


Figure 14.9: The algorithm of Radicchi *et al.* The algorithm of Radicchi *et al.* uses a different measure to identify between-group edges, looking for the edges that belong to the fewest short loops. In many networks, edges within groups belong to many short loops, such as the loops of length three and four labeled “1” and “2.” But edges between groups, such as the edge labeled “3” here, usually do not belong to such loops, because to do so would require there to be a return path along another between-group edge, of which there are, by definition, few.

short loops. Radicchi *et al.* found that loops of length three and four gave the best results. By repeatedly removing edges that belong to small numbers of such loops they were able to accurately uncover communities in a number of example networks.

An attractive feature of this method is its speed. The calculation of the number of short loops to which an edge belongs is a local calculation and can be performed for all edges in time that goes like the total size of the network. Thus, in the worst case, the running time of the complete algorithm that removes all edges one by one will only go as $O(n^2)$ on a sparse network, which is one order of system size faster than the full betweenness-based algorithm and as fast as several of the other methods described in this chapter (though not as fast as the very best of them, such as the spectral method of Section 14.2.3 or the Louvain algorithm of Section 14.2.5).

On the other hand, the algorithm of Radicchi *et al.* has the disadvantage that it only works on networks that have a significant number of short loops in the first place. This restricts the method primarily to social networks, which indeed have large numbers of such loops (see Section 7.3). Other types of networks, such as technological and biological networks, tend to have smaller numbers of short loops, which makes it harder to distinguish between-group edges from within-group ones.

14.5.2 HIERARCHICAL CLUSTERING

The betweenness algorithm of Section 14.5.1 differs from the other community detection algorithms we have considered in producing a hierarchical decomposition of a network into a set of nested communities, as in Fig. 14.8, rather than just a single community division. In this section we look at another algorithm that produces a hierarchical decomposition, one of the oldest of community detection methods, the method of *hierarchical clustering*.¹¹

Hierarchical clustering is not so much a single algorithm as a class of algorithms, with many variations and alternatives. Hierarchical clustering is an agglomerative technique (similar in this respect to the Louvain algorithm of Section 14.2.5), in which we start with the individual nodes of a network and join them together to form groups. The basic idea is to define a measure of similarity or connection strength between nodes, based on the network structure, and then join together the most similar nodes to form groups. We discussed measures of node similarity at some length in Section 7.6 and any of the measures introduced there would be suitable as a starting point for hierarchical clustering, including the cosine similarity of Eq. (7.36), correlation coefficients between rows of the adjacency matrix (Eq. (7.39)), or the Hamming distance of Eq. (7.40). The regular equivalence measures of Section 7.6.2 might also be good choices, although we are not aware of them having been used in this context.

That there are many choices of similarity measures is both a strength and a weakness of the hierarchical clustering method. It gives the method flexibility and allows it to be tailored to specific problems, but it also means the method gives different answers depending on the choices we make, and in many cases there is no way to know in advance if one choice is more correct or will yield more useful information than another. Most often the choice of a similarity measure is determined more by experience or experiment than by argument from first principles.

Once a similarity measure is chosen, we calculate it for all pairs of nodes in the network, then group together those nodes having the highest similarities. The basic strategy for doing this is to start by joining together the pairs of nodes with the highest similarities, forming groups of size two. Then we further join together the groups that are most similar to form larger groups, and so on.

¹¹The word “clustering” is used here as another name for community detection. We have mostly stayed away from using this word in this sense, to avoid confusion with the other sense introduced in Section 7.3, but the name “hierarchical clustering” is a well-established and traditional one, and we use it here in deference to convention.

This, however, brings up another problem: in order to join together the most similar groups we need a measure of group similarity, but what we have is a measure of node similarity. The usual solution to this problem is to combine the node similarities somehow to create similarities for the groups. There are three common ways of doing this, called single-, complete-, and average-linkage clustering.

Consider two groups of nodes, group 1 and group 2, containing n_1 and n_2 nodes respectively. There are n_1n_2 pairs of nodes such that one node is in group 1 and the other in group 2. In the *single-linkage clustering* method, the similarity between the two groups is defined to be the similarity of the *most similar* of these n_1n_2 pairs of nodes. Thus, if the values of the similarities of the node pairs range from 1 to 100, the similarity of the two groups is 100. This is a very lenient definition of similarity: only a single node pair need have high similarity for the groups to be considered similar. (This is the origin of the name “single-linkage clustering”—similarity between groups is a function of the similarity between only the single most similar pair of nodes.)

At the other extreme, *complete-linkage clustering* defines the similarity between two groups to be the similarity of the *least similar* pair of nodes. If the similarities of the nodes range from 1 to 100, then the similarity of the groups is 1. By contrast with single-linkage clustering this is a very stringent definition of group similarity: every single node pair must have high similarity for the groups to have high similarity (hence the name “complete-linkage clustering”).

In between these two extremes lies *average-linkage clustering*, in which the similarity of the two groups is defined to be the mean similarity among all the pairs of nodes. Average-linkage clustering is probably the most satisfactory choice of the three, being a moderate one—not extreme in either direction—and depending on the similarity of all node pairs and not just of the most or least similar pair. It is, however, relatively rarely used, for reasons that are not entirely clear.

The full hierarchical clustering method is as follows:

1. Choose a similarity measure and evaluate it for all node pairs.
2. Assign each node to a group of its own, consisting of just that one node. The similarities of these groups are just the similarities of the nodes.
3. Find the pair of groups with the highest similarity and join them together into a single group.
4. Calculate the similarity between the new composite group and all others using one of the three methods above (single-, complete-, or average-linkage clustering).
5. Repeat from step 3 until all nodes have been joined into a single group.

How fast is this algorithm? The most demanding part of the algorithm is the calculation of the new similarities.¹² Let us consider the three cases separately. For single-linkage clustering the similarity of two groups is equal to the similarity of their most similar pair of nodes. In this case, when we join groups 1 and 2 together, the similarity of the composite group to another group 3, is the greater of the similarities of 1 with 3 and 2 with 3, which can be found in $O(1)$ time.

For complete-linkage clustering the similarity of the composite group is the lesser of the similarities of 1 with 3 and 2 with 3, which can also be found in $O(1)$ time.

The average-linkage case is only slightly more complicated. Suppose as before that the groups 1 and 2 to be joined have n_1 and n_2 nodes, respectively. Then if the similarities of 1 with 3 and 2 with 3 were previously σ_{13} and σ_{23} , the similarity of the composite group with another group 3 is given by the weighted average

$$\sigma_{12,3} = \frac{n_1\sigma_{13} + n_2\sigma_{23}}{n_1 + n_2}. \quad (14.51)$$

Again this can be calculated in $O(1)$ time.

On each step of the algorithm we must calculate similarities in this way for the composite group with every other group, of which there are $O(n)$ in the worst case. Hence the recalculation of similarities will take $O(n)$ time. A naive search through the similarities to find the greatest one, on the other hand, takes time $O(n^2)$, since there are $O(n^2)$ pairs of groups to check, so at first sight it appears that this—and not the recalculation—will be the most time-consuming step in the algorithm. We can speed things up, however, by storing the similarities in a heap, a data structure that allows us to add and remove entries in time $O(\log n)$ and find the greatest one in time $O(1)$. This slows the recalculation of the similarities to $O(n \log n)$ but speeds the search for the largest one to $O(1)$.

The whole process of joining groups has to be repeated $n - 1$ times until all nodes have been joined into a single group. (To see this, simply consider that the number of groups goes down by one every time two groups are joined, so it takes $n - 1$ joins to go from n initial groups to just a single one at the end.) Thus the total running time of the algorithm is $O(n^3)$ in the naive implementation or

¹²This assumes that the initial similarity values for all pairs of nodes can be calculated relatively quickly—for instance in time $O(n^2)$. If this is not possible, if the initial calculation takes much longer, then it is this initial calculation that will determine the overall running time of the algorithm and not the hierarchical clustering procedure itself.

We encountered heaps previously in Section 8.6 in our discussion of Dijkstra's algorithm. An introduction to their properties and working can be found in, for instance, Cormen *et al.* [122].

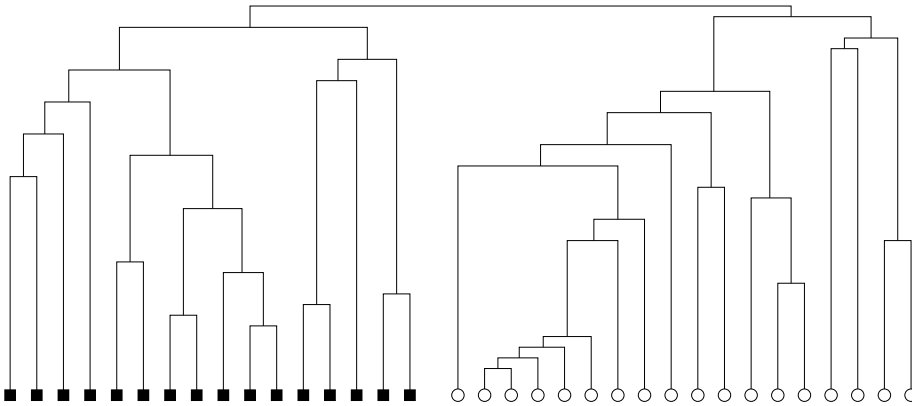


Figure 14.10: Partitioning of the karate club network by average linkage hierarchical clustering. This dendrogram is the result of applying the hierarchical clustering method described in the text to the karate club network of Fig. 14.3, using cosine similarity as our measure of node similarity. The shapes and colors of the nodes represent the two known factions in the network, as in Fig. 14.3.

$O(n^2 \log n)$ if we use a heap.¹³

And how well does it work in practice? The answer depends on which similarity measure one chooses and which linkage method, but a typical application, to the karate club network examined previously in Section 14.2.2, is shown in Fig. 14.10. This figure shows what happens when we apply average-linkage clustering to the karate network using cosine similarity as our similarity measure. The figure shows the dendrogram that results from such a calculation and we see that there is a clear division of the dendrogram into two communities that correspond perfectly to the two known groups in the network.

Hierarchical clustering does not always work as well as this. In particular, though it is often good at picking out the cores of groups, where the nodes are strongly similar to one another, it tends to be less good at assigning peripheral nodes to appropriate groups. Such nodes may not be strongly similar to any others and so tend to get left out of the clustering process until the very end. A common outcome of hierarchical clustering is thus a set of tightly knit cores

¹³For the special case of single-linkage clustering, there is a slightly faster way to implement the algorithm that makes use of a so-called union/find technique and runs in time $O(n^2)$. In practice the performance difference is not very large but the union/find method is considerably simpler to program. It is perhaps for this reason that single-linkage is more often used than complete- or average-linkage clustering.

surrounded by single nodes or smaller groups, though such a decomposition may still contain a lot of valuable information about the network structure.

We have here seen just a few of the best known methods for performing community detection. There are numerous others in the large and growing literature on this topic. We must turn our attention to other things for the remainder of this chapter, but if you are interested in learning more about community detection the review articles by Fortunato [183] and Fortunato and Hric [185] provide useful overviews.

14.6 MEASURING ALGORITHM PERFORMANCE

As we have seen, there are many different approaches and algorithms for community detection. We have examined half a dozen of the most widely used in this chapter and there are many others to be found in the research literature. A natural question that arises, therefore, is which algorithm is best.

There is no single answer to this question since algorithms differ not only in the quality of their results but also in their running speed, their mathematical rigor, and the exact type of results they return. If one wants guarantees of good performance, for instance, then an approach like the maximum likelihood method of Section 14.4.1 might be attractive. But if one wishes to analyze very large networks, then one might be willing to sacrifice some rigor in return for a faster algorithm like the Louvain method of Section 14.2.5.

The primary criterion for choosing between different methods, however, must surely be the quality of the results they return, so our first question in evaluating any method is how well it does at actually finding communities. If its performance is poor, then it is unlikely to be anyone's first choice, no matter what other features it might have.

How does one measure the performance of a community detection algorithm? There are two basic approaches. The first is to test the algorithm on real-world networks whose community structure is widely agreed upon. The second is to test on artificially generated networks with specific known community structure planted within them. Both approaches have their advantages. The first has the benefit of using real-world networks, which present a more realistic test of algorithm performance. On the other hand, it is often hard to know exactly what the true community structure is for real-world networks. In artificial networks, by contrast, we typically know exactly where the communities are, because we put them there. On the other hand, artificial networks are inevitably less realistic, and hence give a less authentic test of real-world algorithm performance. In practice, both methods are widely used and they

are to some extent complementary. In this section we give examples of each of them.

14.6.1 TESTS ON REAL-WORLD NETWORKS

Tests of community detection on real-world networks rely on being able to find example networks where we know—or believe we know—the true division into communities, sometimes called the *ground truth* division. In most cases, the ground truth is established through a combination of insider knowledge of the network in question and consensus results from the application of many different community detection methods to the same network. A classic example in this regard—indeed *the* classic example—is one we’ve already seen, the “karate club” network of Fig. 14.3.

Karate club network: As discussed in Section 14.2.2, this network represents the pattern of friendships among a group of students at a university karate club. The friendships were recorded by a researcher, Wayne Zachary, as part of an anthropological study [479]. As luck would have it, a dispute arose within the club during the study and the club split in two. It is the two factions in this split, as reported by Zachary, that form the ground truth for the network, bolstered by the fact that these same factions have been reproduced in repeated analyses of the network using a broad range of community detection algorithms. So certain are we by now of the true communities in the karate club network that the network has become a kind of qualifying exam for any new detection algorithm. A T-shirt popular among network science researchers carries a picture of the karate club network accompanied by the caption, “If your method doesn’t work on this network then go home.”¹⁴

Several other networks have also found their way into network researchers’ collective consciousness as standardized tests for community detection.

Dolphin social network: An example analogous in many ways to the karate club is the dolphin social network of Lusseau *et al.* [316]. As discussed in Section 4.3, many animal species form persistent social networks and dolphins are a prime example: numerous studies have shown that dolphins form long-lasting “friendships,” meaning in practice that the same pair of dolphins will be seen together often. Lusseau *et al.* spent some years watching a particular group

¹⁴There also exists a tongue-in-cheek society, the “Karate Club Club,” membership in which is bestowed upon the first speaker to use the karate club network as an example in their presentation during any conference attended by the previous recipient of the same honor. A small plastic trophy is passed from hand to hand to commemorate the occasion. The present author has the dubious distinction of being the fifth holder of the trophy.

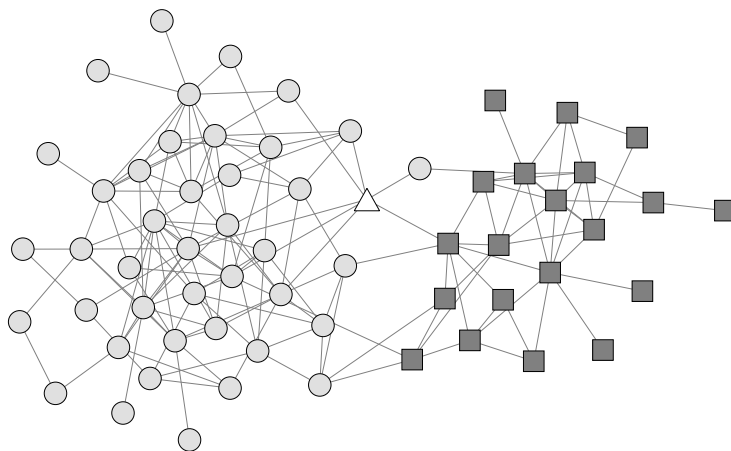


Figure 14.11: Social network of dolphins. A network of 62 bottlenose dolphins in Doubtful Sound, New Zealand, compiled from observations by Lusseau *et al.* [316]. The edges connect pairs of dolphins that were observed together frequently. During the course of the study, one dolphin, represented by the triangle in the center, disappeared, after which the remaining dolphins split into two separate groups, indicated by the circles and squares.

of 62 dolphins off the coast of New Zealand, learning to recognize them by their markings and recording observed pairings. Based on these observations, they assembled the social network shown in Fig. 14.11. Part of the way through the study, however, the group split in two, apparently because of the departure of a pivotal member of the group, and, as with the karate club, the membership of the two factions was recorded and serves as the ground truth for tests of community detection on the network. Repeated studies have found that this ground truth can be duplicated, at least approximately, by community detection algorithms acting on the complete network as measured before the split occurred.

Political blogs: A more challenging test of community detection can be found in the widely studied political blog network of Adamic and Glance [4]. This is a web network in which the nodes represent 1494 weblogs, or blogs for short, about US politics, as observed in the run-up to the 2004 US presidential election. Edges in the network represent hyperlinks between blogs and are in principle directed, although in most community detection studies the directions are ignored and the network treated as undirected. Most studies have also been

See Section 3.1 for a discussion of web networks.

confined to the largest (weakly) connected component of the network, which contains 1225 nodes.

It is widely observed that like-minded political websites, those on the same side of the political aisle, link to one another more often than those on opposite sides, so the blog network should show conventional community structure with respect to political orientation. In the US political system there are two major political parties, Democrats on the left and Republicans on the right, so we expect to see two communities. Adamic and Glance determined the political outlook of each of the blogs in their data set by hand, looking at the text of blog posts to determine which party each blog supported. (Every blog was designated as being either left or right; none were considered centrist.) These designations form the ground truth for the network and numerous studies have found that community detection algorithms can, with good though not perfect accuracy, recover this ground truth from the network structure. The network provides a greater challenge for community detection both by being significantly larger than the karate club and dolphin examples discussed above and by having a highly skewed degree distribution—like most web networks, this one has a long-tailed degree distribution that roughly follows a power law.

American football competition: The three previous networks all have just two communities, which is the simplest possible situation (since there is no community detection to be done in a network with only one community). It is important, however, to test algorithms in more complex situations too. The widely studied “American college football” network provides such a test. This network represents the pattern of games between 115 top-ranked university teams in the annual NCAA American football competition during the fall of 2000. (One could construct a network for any year, but the network for the year 2000 has by consensus become the standard test.)

University-level competition in American football is (to the surprise of some non-Americans) big business in the US. Games are often nationally televised and the best team coaches command salaries of millions of dollars. There are enough highly competitive teams that they cannot all play each other in a single season, so competition has traditionally been divided into smaller groups called “conferences,” typically having about a dozen members and organized roughly along geographic lines. In any given season, the majority of games are between teams in the same conference, although a small fraction are interconference games.

In the network representation of this process the nodes correspond to teams, the edges correspond to games between teams, and the conferences provide the ground-truth communities. The conferences and their membership vary some-

See Section 10.3 for a discussion of the degree distribution of the World Wide Web.

what from year to year, but in the 2000 season there were eleven conferences, plus a small number of independent teams that belonged to no conference. Repeated studies have shown that it is possible to identify the membership of the conferences accurately (and even to pick out the independents) by applying community detection methods.

14.6.2 ARTIFICIAL TEST NETWORKS

While real-world networks provide a test of community detection performance in life-like situations, they are limited in that there are a relatively small number of accepted test networks, their ground-truth community structure is not always 100% certain, and their structure cannot be varied to probe algorithm behavior. We cannot, for instance, make their communities larger or smaller, or change their number.

A solution to these problems is to use artificial test networks, sometimes called *synthetic networks*, which contain a specific level or type of community structure planted within them. We can then test to see whether our community detection algorithms can detect this planted structure. While artificial networks are generally less realistic than their real-world counterparts, their structure can be varied in any way we please and we can generate as many of them as we want, which gives us considerable flexibility to quantify the performance of our algorithms. In practice, algorithms are often tested on both real and synthetic networks, to make the most of the advantages each has to offer.

There are a number of ways one might envisage generating artificial networks with embedded community structure, but in practice there are two main approaches that find common use.

Stochastic block models: The most common approach to generating artificial test networks for community detection is to use a stochastic block model of the kind introduced in Section 12.11.6. Recall that in the standard stochastic block model one divides n nodes into some number q of groups (not necessarily of equal size, although in practice most tests do use equal sizes), then places an edge between each pair of nodes with independent probability p_{rs} , where r and s are the groups to which the nodes belong. The probabilities p_{rs} form a $q \times q$ matrix of parameters whose values determine the community structure. In the most common variant of the model the parameters take just two values

$$p_{rs} = \begin{cases} p_{\text{in}} & \text{if } r = s, \\ p_{\text{out}} & \text{if } r \neq s. \end{cases} \quad (14.52)$$

If $p_{\text{in}} > p_{\text{out}}$ then edges are more likely within groups than between them and the network has traditional community structure with dense groups of nodes

connected by sparser inter-group edges.

Changing the values of p_{in} and p_{out} allows us to vary the difficulty of the community detection problem. If p_{in} is much larger than p_{out} then the communities in the network should stand out clearly and most algorithms should have little difficulty detecting them. However, if the difference between the two probabilities is small then there will be little to distinguish in-group edges from between-group ones and the community structure will be hard to pick out. Indeed if we let the difference go to zero, so that $p_{\text{in}} = p_{\text{out}}$, then all edges in the network have equal probability and by definition there is no community structure at all. In this limit, therefore, all algorithms must fail.

Several other models commonly used for generating test networks are also effectively block models, even if they are not given that name. Examples include the *planted partition model* of Condon and Karp [120], which is widely used in computer science, and the “four groups” test of Girvan and Newman [204], which is popular in the physics literature.

The LFR benchmark: The stochastic block model is simple to describe and simple to use, but it does not present a very realistic challenge for community detection algorithms. The networks it produces have a Poisson degree distribution within each group and the groups are, in most cases, of equal sizes. Real-world networks, on the other hand, typically have strongly right-skewed degree distributions (see Section 10.3) and unequally sized groups. One can fix both of these issues easily enough: one can introduce a different degree distribution, for example, by using the degree-corrected stochastic block model (Section 12.11.6), and there is nothing other than tradition stopping us from choosing unequal group sizes. In recent years it has become common to test community detection algorithms against more realistic synthetic networks that remedy the shortcomings of the stochastic block model. For historical reasons, however, tests are not usually performed using the degree-corrected block model but a slight variant, called the *LFR model* or *LFR benchmark*, after its inventors, Lancichinetti, Fortunato, and Radicchi [286,287].

In this model one again divides n nodes among some number q of groups, but the groups are now of varying sizes. Motivated by observations of communities in certain real-world networks [110,223], the sizes of the groups are chosen from a power-law distribution, subject only to the constraint that they must sum to n . Once nodes have been assigned to groups, edges are placed between them, but they are not placed independently. Instead one first chooses the degree of each node from another power-law distribution (motivated again by the common appearance of such power laws in real degree distributions—see Section 10.3). One also chooses a parameter value μ , which measures the

See Section 10.4 for an introduction to power-law distributions.

fraction of a node's edges that connect to nodes in different groups. Then for each node having degree k , one gives it μk connections to other groups (rounded to the nearest integer) and $(1 - \mu)k$ connections to its own group. The connections themselves are constructed in a manner similar to the configuration model of Section 12.1: one can think of each node's connections as "stubs" of edges attached to the node, and one draws stubs at random in pairs and connects them to form edges, while being careful to match in-group stubs with others in the same group, and between-group stubs with others in different groups.

The end result is an artificial network with a power-law degree distribution and power-law distribution of community sizes, which poses a more demanding test for community detection than the simpler and more uniform stochastic block model.

14.6.3 QUANTIFYING PERFORMANCE

Given a particular test network, how can we say if an algorithm is doing a good job of finding the known community structure in that network? If the structure found by the algorithm is identical to the agreed ground truth, then we can pat ourselves on the back and claim victory, but in most cases we will not be so lucky. A more common outcome is that an algorithm finds communities that look somewhat similar to the ground truth, but are not exactly identical. We would like a way to quantify how close to the correct answer such a result comes. There are three commonly used methods for doing this.

Fraction of correctly classified nodes: The simplest measure of success is just to count the fraction of nodes that are classified into the correct groups. Suppose our ground truth tells us that nodes 1 to 5 of a 10-node network are in group 1 and the remainder, nodes 6 to 10, are in group 2, so that the vector \mathbf{g} consisting of n group assignment variables g_i looks like

$$\mathbf{g} = (1, 1, 1, 1, 1, 2, 2, 2, 2, 2). \quad (14.53)$$

Now we run our favorite community detection algorithm on the network. The algorithm does a pretty good job, finding a division of the network that is close to the ground truth, but it gets a couple of nodes wrong. Its output look likes this:

$$\mathbf{g} = (1, 2, 1, 1, 1, 2, 2, 1, 2, 2). \quad (14.54)$$

In this case the algorithm gets eight out of ten nodes right, so the fraction of correctly classified nodes is 0.8.

This approach works reasonably well, but there is a catch. Sometimes when we run our community detection algorithm we might get a result like this:

$$\mathbf{g} = (2, 1, 2, 2, 2, 1, 1, 2, 1, 1). \quad (14.55)$$

This result corresponds to the exact same division of the network as (14.54). The only thing different is that the group labels 1 and 2 are reversed. The group labels themselves are meaningless—which group we call number 1 and which we call number 2 is entirely arbitrary. All we actually care about is the group division. In this sense Eq. (14.55) is every bit as good as Eq. (14.54). If we compare (14.55) with the ground truth in (14.53), however, we find that only two nodes are classified correctly, so the fraction of correct nodes is 0.2.

To get around this problem, one commonly calculates the fraction of correctly classified nodes as the maximum over all permutations of the group labels. That is, one considers each possible way of labeling the groups found by the algorithm and compares each one to the ground truth, choosing from among the possibilities the one that gives the largest number of correctly classified nodes. Alternatively, we can if we prefer permute the labels in the ground truth—either approach gives the same result.

It might also be the case that the number of groups found by the algorithm is not equal to the number of groups in the ground truth. The same permutation approach works in this case, except that one must always permute the larger of the two sets of labels—the output of the algorithm or the ground truth. In the example above, for instance, if our algorithm found three groups with

$$\mathbf{g} = (2, 3, 2, 2, 2, 1, 1, 3, 1, 1), \quad (14.56)$$

then the right way to calculate the fraction of correctly classified nodes would be to maximize the fraction over all ways of labeling the three groups, which gives 0.8 again in this case.

Note that in the process of permuting the labels of q groups, every node by definition gets classified correctly in a fraction $1/q$ of the permutations, so the average over all permutations of the fraction of correctly classified nodes is always $1/q$. That in turn means that the maximum over all permutations can never be less than $1/q$, since the maximum of a set of numbers is never less than their mean. Thus the score for, say, a two-group division can never be less than $\frac{1}{2}$ and our results should be evaluated in this light. For a two-group division, a score of 0.6—meaning that 60% of nodes are classified correctly—might at first sight appear promising, but in fact it is only a little better than the minimum possible value of 0.5.

Another way to think about this is that an entirely random assignment of nodes to groups would on average get $1/q$ of the nodes right. So by comparing

our results to a baseline value of $1/q$ we are asking how much better our community detection algorithm does than a random roll of the dice.

Rand index: The need to maximize the fraction of correctly classified nodes over permutations of the group labels is an annoying (and sometimes time-consuming) complication. An alternative measure of performance that avoids this issue is the *Rand index* [399].

The Rand index relies on the observation that two nodes that are placed in the same group will remain in the same group no matter how we permute the group labels, and likewise for nodes in different groups. The Rand index measures how often nodes that are in the same (or different) groups in the ground truth are also assigned to the same (or different) groups by our community detection algorithm.

In a network of n nodes there are $\binom{n}{2}$ pairs of distinct nodes. Among all of those pairs, let s be the number that are in the same group in the ground truth and are also (correctly) placed in the same group by our community detection algorithm. If we once again denote the community to which node i is assigned by g_i , and if we denote the ground-truth community for the same node by t_i , then

$$s = \sum_{i < j} \delta_{g_i g_j} \delta_{t_i t_j}, \quad (14.57)$$

where by summing over node pairs with $i < j$ we ensure that each pair is counted only once and that pairs with $i = j$ are excluded.

Similarly, let d be the number of pairs of nodes that are in different groups in the ground truth and are also placed in different groups by the algorithm:

$$d = \sum_{i < j} (1 - \delta_{g_i g_j})(1 - \delta_{t_i t_j}). \quad (14.58)$$

Thus $s + d$ is the total number of pairs correctly placed in the same or different groups by the algorithm. The Rand index is this number expressed as a fraction of the total number of pairs:

$$R = \frac{s + d}{\binom{n}{2}}. \quad (14.59)$$

Values of the Rand index lie in the range $0 \leq R \leq 1$, with high values indicating that the algorithm has assigned a large fraction of node pairs correctly to either the same or different groups, i.e., that it has accurately detected the community structure.

The Rand index has the advantage that it does not depend on the labels used to identify the groups, making it simple to calculate. There is no need to maximize over permutations of the labels. On the other hand, the Rand index

is somewhat harder to interpret than the fraction of correctly classified nodes. It is perhaps for this reason that it has, so far, found relatively little use in the community detection literature.

Normalized mutual information: A third—and widely used—measure for comparing the output of a community detection algorithm to ground truth is the *normalized mutual information*, a measure based on ideas from information theory. Information theory, in its simplest form, is a way of quantifying the information content of strings of letters, numbers, or other symbols. We can consider a vector \mathbf{g} of group assignments, as output by one of our community detection algorithms, as just such a string of symbols.

Specifically, let us again denote the group to which node i is assigned by our community detection algorithm by g_i and the ground-truth community by t_i , and let $P(t|g)$ be the probability that the true, ground-truth community for a node is t , given that it was assigned to group g by our community detection algorithm. We can estimate this probability simply by going through all nodes assigned to group g and finding the fraction with ground truth t :

$$P(t|g) = \frac{\sum_i \delta_{t_i t} \delta_{g_i g}}{\sum_i \delta_{g_i g}}. \quad (14.60)$$

The *conditional entropy* of the complete vector of ground-truth assignments \mathbf{t} , given the group assignments \mathbf{g} found by the algorithm, is defined to be

$$H(\mathbf{t}|\mathbf{g}) = - \sum_g P(g) \sum_t P(t|g) \log P(t|g), \quad (14.61)$$

where $P(g)$ is the probability that a node is assigned to group g by the algorithm, which we can estimate from $P(g) = \sum_i \delta_{g_i g} / n$.

The conditional entropy tells us the amount of additional information contained in the ground truth, if we already know the group assignments found by the algorithm. In other words, if we run our algorithm and get an assignment \mathbf{g} , and then someone gives us the true assignments of the nodes, the conditional entropy tells us how much more we learn about group membership from those true assignments beyond what we already got from our algorithm.

For example, if the ground truth and the output of the algorithm are precisely identical—if the algorithm has performed perfectly—then all probabilities $P(t|g)$ are either one or zero, which means that $P(t|g) \log P(t|g)$ is zero¹⁵ and hence the conditional entropy of Eq. (14.61) is also zero. This is the correct

¹⁵Strictly $x \log x$ is undefined when $x = 0$, but the limit as $x \rightarrow 0$ is well defined and is equal to zero.

We encountered some results from information theory previously in Section 14.4.1.

answer: in this case the algorithm tells us everything, so the ground truth gives us zero additional information.

Conversely, suppose our algorithm fails completely, returning a group assignment that is totally unrelated to the true assignment. In that case $P(t|g) = P(t)$, independent of g , and

$$H(\mathbf{t}|\mathbf{g}) = - \sum_g P(g) \sum_t P(t) \log P(t) = - \sum_t P(t) \log P(t), \quad (14.62)$$

where we have made use of $\sum_g P(g) = 1$. The quantity

$$H(\mathbf{t}) = - \sum_t P(t) \log P(t) \quad (14.63)$$

is the (unconditional) entropy of the ground truth, the total information contained in the ground truth. In other words, when the algorithm fails, the amount of additional information we get from the ground truth is equal to the total amount of information in the entire group assignment—we get a maximal amount of information from the ground truth, since we learn nothing from the algorithm and the ground truth tells us everything.

Note, crucially, that these results do not depend on how the groups are labeled, either by the algorithm or in the ground truth. In particular, the same group does not have to be given the same label in both for us to correctly calculate the conditional entropy. This is one of the key advantages of this approach. It means that one does not have to search through all possible permutations of the group labels to find the correct one. Any permutation will do.

Thus the conditional entropy provides an information theoretic measure of how well our community detection algorithm performs, ranging between its lowest and highest values as the performance of the algorithm goes from perfect to hopeless.

It's not a very convenient measure, however, since it is small when the algorithm does well and large when it does poorly. To fix this, we can flip the measure around by subtracting it from its maximum value, Eq. (14.63), which gives the quantity known as the *mutual information*:

$$I(\mathbf{t}; \mathbf{g}) = H(\mathbf{t}) - H(\mathbf{t}|\mathbf{g}). \quad (14.64)$$

This quantity is now zero when the algorithm totally fails and takes its maximal value of $H(\mathbf{t})$ when performance is perfect.

This is still not ideal, however, since the scale on which the algorithm's success or failure is measured runs from zero to $H(\mathbf{t})$, which will vary from

network to network. Usually, therefore, we normalize the measure so that its value falls always between zero and one. You might imagine we would do this by simply dividing by $H(\mathbf{t})$, but commonly we make a slightly different choice. It is straightforward to show that $I(\mathbf{t}; \mathbf{g})$ is symmetric with respect to its arguments, meaning that¹⁶

$$I(\mathbf{t}; \mathbf{g}) = I(\mathbf{g}; \mathbf{t}) = H(\mathbf{g}) - H(\mathbf{g}|\mathbf{t}), \quad (14.65)$$

and hence it is bounded above not only by $H(\mathbf{t})$ but also by the entropy $H(\mathbf{g})$ of the assignment found by the algorithm. Thus it can have a value that is at most equal to the smaller of the two entropies and we can safely normalize by the minimum of the two and be guaranteed a result that falls in the interval from zero to one:

$$N(\mathbf{t}; \mathbf{g}) = \frac{I(\mathbf{t}; \mathbf{g})}{\min[H(\mathbf{t}), H(\mathbf{g})]}. \quad (14.66)$$

This is one version of the normalized mutual information, though it's not the version most commonly used. The standard version, proposed by Danon *et al.* [128], normalizes by the average $\frac{1}{2}[H(\mathbf{t}) + H(\mathbf{g})]$ of the entropies (rather than the minimum), giving

$$N(\mathbf{t}; \mathbf{g}) = \frac{2I(\mathbf{t}; \mathbf{g})}{H(\mathbf{t}) + H(\mathbf{g})}. \quad (14.67)$$

This is the form in which the normalized mutual information is most commonly used. It has two nice features:

1. It runs from a value of zero when the algorithm fails to a value of one when the algorithm works perfectly. To see that the latter is true, note then when the algorithm works perfectly we have $\mathbf{g} = \mathbf{t}$, $H(\mathbf{g}) = H(\mathbf{t})$, and

¹⁶The proof looks like this:

$$\begin{aligned} I(\mathbf{x}; \mathbf{y}) &= - \sum_x P(x) \log P(x) + \sum_y P(y) \sum_x P(x|y) \log P(x|y) \\ &= - \sum_{xy} P(x, y) \log P(x) + \sum_{xy} P(x, y) \log \frac{P(x, y)}{P(y)} \\ &= - \sum_{xy} P(x, y) \log P(y) + \sum_{xy} P(x, y) \log \frac{P(x, y)}{P(x)} \\ &= - \sum_y P(y) \log P(y) + \sum_x P(x) \sum_y P(y|x) \log P(y|x) = I(\mathbf{y}; \mathbf{x}). \end{aligned}$$

$H(\mathbf{t}|\mathbf{g}) = 0$, so that

$$N(\mathbf{t}; \mathbf{g}) = \frac{2I(\mathbf{t}; \mathbf{g})}{H(\mathbf{t}) + H(\mathbf{g})} = \frac{H(\mathbf{t}) - H(\mathbf{t}|\mathbf{g})}{H(\mathbf{t})} = 1. \quad (14.68)$$

2. It maintains the symmetry of the mutual information with respect to its two arguments, which would be lost if we simply normalized by a factor of $H(\mathbf{t})$.

This choice of normalization for the mutual information is somewhat arbitrary. Had we normalized by $\min[H(\mathbf{t}), H(\mathbf{g})]$ as in Eq. (14.66) then the two properties above would still be satisfied. Some authors even normalize by the *maximum* of the two entropies, which also has the same properties.

See McDaid *et al.* [327] for a discussion of the several ways in which the mutual information can be normalized.

The form in Eq. (14.67) is, however, the most widely accepted one and is commonly used as a measure of the performance of community detection algorithms.

14.6.4 COMPARISON OF COMMUNITY DETECTION ALGORITHMS

There have been a number of studies comparing the performance of different community detection methods using the techniques described in the preceding sections [128, 286, 376, 476]. There is no universal consensus about which algorithm comes out on top, but overall the modularity maximization and InfoMap methods seem to get the most nodes. The early study of Danon *et al.* [128], which used networks generated with the stochastic block model, and the more recent study of Yang *et al.* [476], which used the LFR benchmark, both find the best performance for modularity maximization methods (Section 14.2), although with different specific techniques for performing the maximization. (Danon *et al.* find in favor of maximization by simulated annealing, while Yang *et al.* favor the Louvain algorithm of Section 14.2.5.) Meanwhile, the study by Orman *et al.* [376], which uses the LFR benchmark again, finds best performance for the InfoMap method, as do Lancichinetti and Fortunato [286], though the latter place the Louvain algorithm a close second. Overall, it appears that either modularity maximization or InfoMap is a solid choice for practical community detection. Both return high-quality results and are fast enough for application to very large networks.

No other algorithm consistently scores well in comparisons. However, none of the studies published so far has included inference methods, as in Section 14.4.1, in their tests. On theoretical grounds one can expect these methods to perform well, and indeed on test networks generated using the stochastic block model they should be unbeatable, since they are formally optimal for this problem [139]. It will be interesting to see the results of

quantitative performance tests on this class of algorithms when they become available.

14.7 DETECTING OTHER KINDS OF NETWORK STRUCTURE

We have so far devoted our attention in this chapter to questions about just one type of structure in networks, namely community structure: what it is, what it means, and how to detect it. However, while community structure is certainly the best-studied form of large-scale network structure, it is by no means the only one. There are many other types of structure that are observed to occur in at least some networks and that can tell us much of interest. In this section we introduce some of these, along with methods for their detection.

14.7.1 OVERLAPPING COMMUNITIES

Not really an entirely new form of structure, but rather a variation on the community structure theme, our first example is *overlapping communities*. In traditional community structure of the type discussed in previous sections and depicted, for instance, in Fig. 14.3 on page 503, every node belongs to exactly one group and one group only. There are situations, however, where it makes sense to allow nodes to belong to more than one group. In a social network of friends, for example, a person might belong to a group of work friends, a group of old college friends, and a group of family members and family friends. All of these might be well-defined groups in the community structure sense, and yet the same person can belong to all of them, meaning that the groups overlap: the common member (or members) of the groups constitute the overlap—see Fig. 14.12.

A range of methods have been proposed for detecting overlapping communities. Perhaps the first was the *CFinder* algorithm of Palla *et al.* [379], which works by finding cliques within a network. *CFinder* performs poorly, however, with networks that have few cliques, of which there are many examples, and it has as a result been largely superseded in recent years by other methods.

An alternative and more general approach is to make use of methods of statistical inference akin to those of Section 14.4. One simply needs to define a suitable random-graph-style model with overlapping communities, then fit it to observed data using an appropriate maximum likelihood method. (Indeed this is a common approach for all of the forms of structure we discuss in this

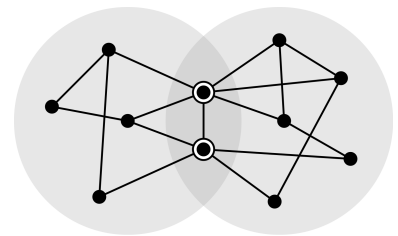


Figure 14.12: Overlapping communities. In this small network there are two communities (shaded circles) and the two nodes in the middle belong to both, so the communities overlap.

A clique is a subset of nodes that are all connected to each other—see Section 7.2.1.

section. One of the advantages of the inference method is that it can be easily adapted to the detection of any kind of structure for which one can imagine a model.)

A number of models for overlapping communities have been proposed. Perhaps the best known is the *mixed-membership stochastic block model* of Airoldi *et al.* [12]. In this model nodes have weights or strengths that measure how strongly they are associated with each possible group: every node i has a set of parameters π_{ir} representing the fraction by which the node belongs to each group r . In a network with three groups, for instance, a node might belong $\frac{1}{2}$ to group 1, $\frac{1}{3}$ to group 2, and $\frac{1}{6}$ to group 3. These fractions generalize the group membership variables of the ordinary stochastic block model, in which every node effectively belongs to one group with weight 1 and all other groups with weight 0. As in the ordinary block model, we also define a set of probabilities p_{rs} for edges between nodes in groups r and s . Then to each pair of nodes i, j we assign groups r, s at random with probabilities π_{ir} and π_{js} respectively, then we place an edge between the nodes with the appropriate probability p_{rs} . Nodes can be assigned to different groups when deciding the placement of different edges: for one of my friends the salient point may be that we work together; for another it might be that we went to the same university. The model does not explicitly include degree correction of the kind discussed in Section 14.4.1 for the standard block model, although it would in principle be straightforward to incorporate it. Airoldi *et al.* describe various different methods for fitting their model to data using maximum likelihood methods akin to those of Section 14.4.

An alternative and conceptually attractive approach to detecting overlapping communities is to cluster the *edges* of a network into groups instead of the nodes [8,166]. Returning to our earlier example of an individual who has work friends, college friends, and family friends, an equivalent—and arguably more elegant—way of representing the situation would be to say that the edges of the network are of several different types: work friendships, college friendships, and family friendships. Thus it is not really the nodes that are grouped in this situation, but the edges. If we can come up with an algorithm to accurately collect the edges into their different classes or types, then we will in the process also have found our overlapping groups. This idea is closely related to the concept of multilayer networks discussed in Section 6.7. In effect, we are saying, the friendship network is a multilayer network consisting of a work layer, a college layer, a family layer, and so on, but we are not told which edges belong to which layers. The task of identifying overlapping communities is essentially equivalent to the task of identifying the layers of the multilayer network.

Again, a number of methods have been proposed for performing this task.

Ahn *et al.* [8] employ a hierarchical clustering technique similar to that of Section 14.5.2 to group edges together, while Evans and Lambiotte [166] propose a method based on random walks. Ball *et al.* [37] propose a method based again on statistical inference, where one defines a model in which each edge is of one of several different types and the total probability of an edge between two nodes is the sum of the probabilities of edges of each type. This allows one to write down an expression for the likelihood of an observed network and from it to infer the type of each edge by a maximum likelihood fit.

14.7.2 HIERARCHICAL COMMUNITIES

So far we have considered network communities as a flat, single-level structure. The network breaks apart into communities, but there is no further structure within those communities. For many networks this is clearly incorrect. Often there can be multiple levels of structure, with communities breaking up into subcommunities, then subsubcommunities, and so on—see Fig. 14.13. A network of acquaintances within a corporation, for instance, might have community structure at the level of departments, teams, individual offices, and so on. We have already seen a couple of methods that, to some extent, reveal this kind of hierarchical division: the betweenness-based method of Section 14.5.1, for example, and the hierarchical clustering method of Section 14.5.2. The output of both of these methods can be represented in the form of a tree or dendrogram like the one in Fig. 14.8 on page 531, which is the most common way of representing hierarchical network structure. The top of the tree represents the entire network, which then splits repeatedly into smaller and smaller subgroups. The process ends at the bottom of the tree when the subgroups become so small that each contains only a single node. The dendrogram in Fig. 14.8 is a binary tree, meaning that when they split, communities always split into exactly two parts. This is the most common choice for dendrograms, in part because this is the way the corresponding algorithms work, but there is no reason in principle why we need restrict ourselves to binary trees. One could certainly have a hierarchical structure in which groups split into varying numbers of subgroups, sometimes two perhaps, but sometimes three or more.

The methods of statistical inference introduced in Section 14.4 can also be applied to the detection of hierarchical structure. One defines an appropriate model of the hierarchy, then fits it to the observed network by maximizing the likelihood. An example of such a model is the *hierarchical random graph* of Clauset *et al.* [109]. In this model, which is depicted in Fig. 14.14, one specifies first a dendrogram, in the form of a binary tree, and then a set of probabilities, one at each internal intersection in the tree. The nodes of the actual network

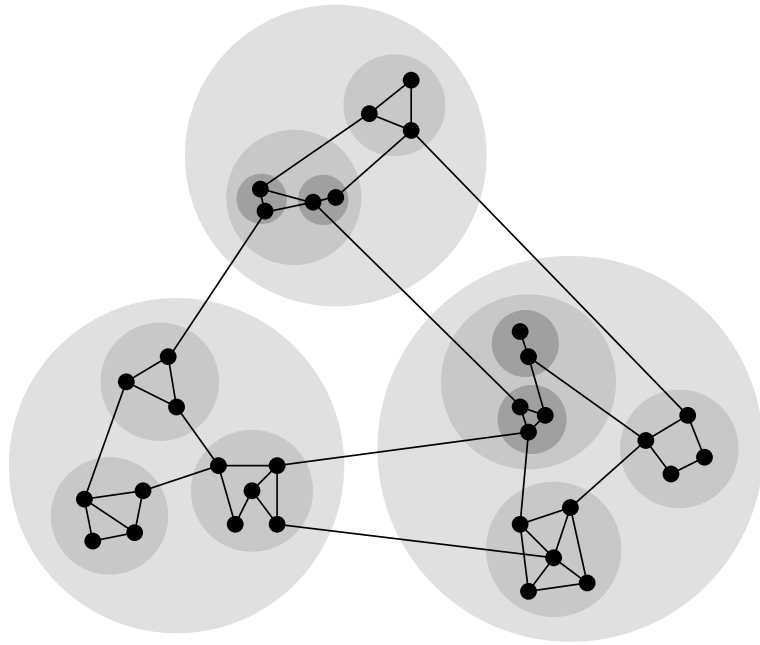


Figure 14.13: A network with a hierarchy of communities within communities. In this sketch of a hierarchical network the network is divided into three large communities, each of which divides into smaller subcommunities, and some of those divide further into subsubcommunities. In principle, this repeated subdivision could continue through a large number of levels.

itself are, as before, represented by the leaves of the tree—the circles at the bottom of the figure—and the probability of an edge between any pair of nodes i, j is equal to the probability stored at the lowest common ancestor of i and j in the tree, i.e., the lowest intersection node that can be reached from both i and j via an upward path through the tree. (An example is shown by the arrows in the figure.) The parameters of the model are the probabilities at the intersections (of which there are $n - 1$ in a network of n nodes) and the tree itself—the structure of the tree parametrizes the hierarchical structure of the network.

Studies indicate that this model appears to fit many empirical network data sets quite well—better than fits to the single-level stochastic block model of Section 12.11.6, or even its degree-corrected variant. The fits accurately reproduce network features like path lengths and clustering coefficients and have been used as the basis for, among other things, “link prediction” algorithms

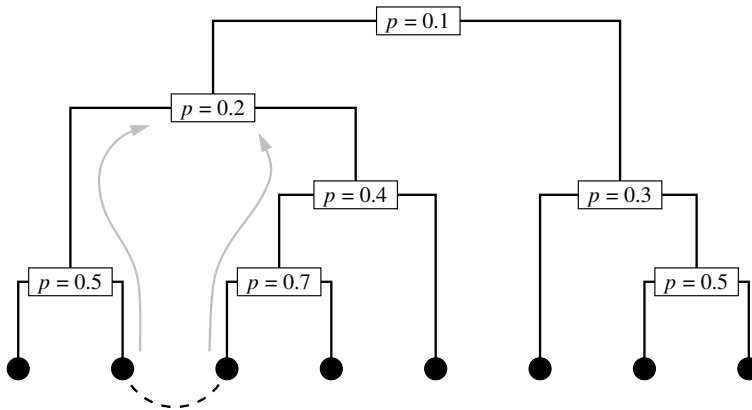


Figure 14.14: The hierarchical random graph. The filled circles at the bottom of this dendrogram represent the nodes of the network being generated, while the dendrogram itself represents the desired hierarchical structure. The probability of an edge between two nodes in the generated network—for instance the edge indicated by the dashed line—is given by the value at the lowest common ancestor of those nodes, which is $p = 0.2$ in this case. The dendrogram itself is never seen in the final network. It and its probabilities exist only to guide the network generation process.

that aim to predict where there might be edges missing from an incompletely observed network [109].

Link prediction was discussed previously in Section 9.4.1.

14.7.3 CORE-PERIPHERY STRUCTURE

Moving away from variants on the community structure idea, *core-periphery structure* [126,414] describes a network whose nodes are divided into a densely connected core and a loosely connected periphery—see Fig. 14.15. One way to think of core-periphery structure is in terms of the average probabilities of edges within and between these two groups of nodes. Suppose we call the two groups group 1 and group 2 and then, by analogy with the stochastic block model of Section 12.11.6, denote the average probabilities of edges within each group by p_{11} and p_{22} and the probability of edges between nodes in different groups by p_{12} . If we have $p_{11} > p_{12} < p_{22}$ then the network has traditional assortative community structure: the probabilities of edges within groups are higher than the probability of edges between groups. Conversely, if $p_{11} < p_{12} > p_{22}$ then we have disassortative structure, with edges within groups being rarer than those between groups (see Section 7.7).

But there is also a third logical possibility: we could have probabilities

$p_{11} > p_{12} > p_{22}$. This is what we mean by core–periphery structure. It is neither assortative nor disassortative in the traditional sense, yet it still has two clear groups, a dense core group (group 1) and a sparse periphery group (group 2), with an intermediate probability of connections between the groups. (The last remaining logical possibility $p_{11} < p_{12} < p_{22}$ is the same thing, just with the group labels 1 and 2 swapped around, so that group 2 is the core and group 1 is the periphery.)

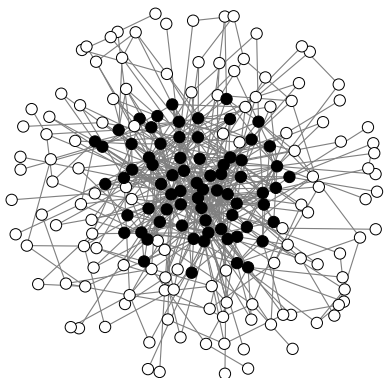


Figure 14.15: Core–periphery structure. A small network with a dense core (solid nodes) and a sparser periphery (open nodes).

In principle, core–periphery structure need not be limited to just two groups. One could have three or more groups ranging from the innermost core with the highest density to the outermost periphery with the lowest, as depicted in Fig. 14.16a. Networks of this kind are sometimes said to have *onion structure* [126], the different groups being analogous to the layers of an onion, from the core outward to the periphery. A network also need not have just one core. It is possible, indeed likely in large networks, for there to be two or more dense core regions that are not directly connected to one another—see Fig. 14.16b.

A very simple method for finding core–periphery structure would be just to assume that the nodes in the core have higher degree than the nodes in the periphery and divide the nodes according to degree. Simple though it is, this method actually works quite well. The results returned by more sophisticated methods typically do not differ that much from this rudimentary degree-based division.

Another simple method is to construct the k -cores of the network (see Section 7.2.2). Recall that a k -core is a group of nodes that each has connections to at least k other members of the group. By definition k -cores form a nested set, with higher k -cores being entirely contained within lower ones, like Russian matryoshka dolls. Moreover, they become denser as k increases, since all nodes within the k th core must have degree at least k . Thus the k -cores provide one possible version of the onion structure described above. The k -cores are also easy to construct—a simple algorithm is given in Section 7.2.2.

Another method for detecting core–periphery structure—of the two-group kind only—has been proposed by Borgatti and Everett [78]. The description in their original paper is rather complicated, but the method boils down to the optimization of a measure somewhat akin to modularity (see Sections 7.7.1 and 14.2). Their basic goal is to find the division of a network into core and periphery that minimizes the number of edges in the periphery. A straight minimization, however, will not work. It is clear that the overall minimum is always achieved by putting all nodes in the core and none at all in the periphery.

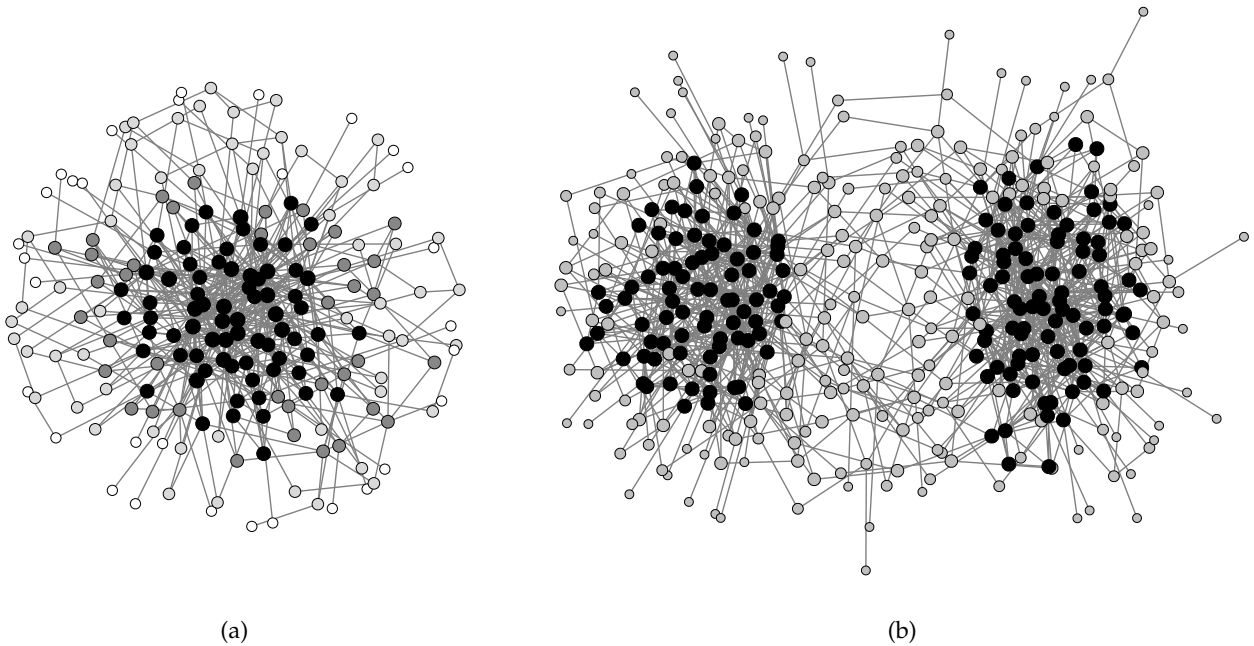


Figure 14.16: Core-periphery structure in two small networks. (a) A network with several layers of core-periphery structure, represented by the varying colors and sizes of the nodes. The darkest color represents the innermost core, which is surrounded by a succession of onion-like layers of lower and lower density. (b) A network with two separate cores that are not directly connected to one another. The cores in both of these examples were found using a k -core decomposition as described in the text.

So, as with modularity, Borgatti and Everett define a score function that is equal to the number of edges in the periphery minus the expected number of such edges if edges were placed at random. If we denote membership of the two groups by $g_i = 0$ when node i is in the core and $g_i = 1$ when i is in the periphery, Borgatti and Everett's score function, which they call ρ , is given by

$$\rho = \frac{1}{2} \sum_{ij} A_{ij} g_i g_j - p \binom{n_p}{2}. \quad (14.69)$$

Here $p = m/\binom{n}{2}$ is the average edge probability if the same number m of edges were placed at random, n_p is the number of nodes in the periphery, and the factor of $\frac{1}{2}$ in front of the sum compensates for double counting of node pairs. If $n_p \gg 1$, which is usually the case, then $\binom{n_p}{2} \approx \frac{1}{2} n_p^2$ and, noting that $n_p = \sum_i g_i$,

we then have

$$\rho = \frac{1}{2} \sum_{ij} (A_{ij} - p) g_i g_j. \quad (14.70)$$

The goal of the method is now to minimize this function, so that the periphery has as few edges as possible in it, relative to what we would expect by chance. Borgatti and Everett performed the minimization using a genetic algorithm, although many other methods could no doubt also be applied.¹⁷

Yet another way to detect core–periphery structure, suggested by our discussion earlier of the average probabilities p_{11} , p_{12} , and p_{22} of connections within and between the core and periphery, is to use an inference method based on fits to a stochastic block model. This approach is essentially the same as community detection using the stochastic block model. You fit an observed network to a block model with two groups and if the best-fit edge probabilities have the form $p_{11} > p_{12} > p_{22}$, then you have found core–periphery structure. In doing this, however, it turns out to be crucial *not* to use the degree-corrected block model that we used for community detection in Section 14.4.1. The reason is that the degree-corrected model effectively factors out any correlation between node degrees and the divisions found by the method. But with core–periphery structure a large part of the information about which nodes are in the core and which are in the periphery resides in the degrees—as we have said, even a simple division by degree does quite a good job. One should therefore use the original, non-degree-corrected stochastic block model for detecting core–periphery structure, so that the fit incorporates degree information into the division it finds [482]. In other respects, however, the method is basically the same as that of Section 14.4.1. The only downside to this approach is that if you fit the model and don’t find that $p_{11} > p_{12} > p_{22}$, then you have failed to find any core–periphery structure. You don’t get to say whether the method finds core–periphery structure or community structure. It just finds whichever gives the highest likelihood.

14.7.4 LATENT SPACES, STRATIFIED NETWORKS, AND RANK STRUCTURE

Another common form of large-scale structure in networks is *latent-space structure*, in which nodes occupy positions along some axis or axes and the presence or absence of edges between nodes depends on the nodes’ positions. In the

¹⁷Borgatti and Everett’s original formulation of ρ was actually the reverse of the one presented here: they maximized the number of edges in all parts of the network except the periphery, i.e., in the core and between the core and the periphery. But since the total number of edges is fixed the two approaches are completely equivalent and the version given here is a little simpler.

most common form of latent-space structure, edges are more likely between nodes that are close together than between those far apart. Networks with this kind of structure are sometimes said to be *stratified*—see the discussion starting on page 206 and the accompanying figure. It is in principle also possible, though rarer, for connections to fall mainly between nodes that are far apart, or for there to be some more complicated dependence of edge probability on distance.

One refers to the axis or axes along which the nodes lie as a latent space, although there is no requirement that it represent actual geometric space. Sometimes it does—people’s homes have geographic positions on the surface of the Earth, for instance, and people are normally more likely to be friends with those who live close to them than with those who live far away. But a latent space can also represent some kind of non-geographic coordinate, such as age for instance—people are also more likely to be friends with others close to them in age.

The idea of latent-space structure is related to the concept of assortative mixing discussed in Section 7.7.2. There we considered the case where nodes possess some characteristic property measured by a scalar variable, which can be thought of as a coordinate in a latent space, and we constructed a measure to quantify the extent to which nodes with similar values tend to be connected by edges. Here we take a contrasting approach and assume we do not know the values on the nodes and ask whether we can guess those values by looking at the network structure.

For instance, the simplest version of this problem asks whether there is any way to arrange the nodes of a given network on a line, such that most edges join nodes that are close together. Such an arrangement of the nodes is called an *embedding* of the network in the one-dimensional space of the line. One can also consider embeddings in spaces of two or more dimensions. One way to formalize the problem of finding a good embedding is to write an expression for a suitable objective function, such as the sum of the squared distances spanned by the network’s edges, and then minimize that function over all choices of node positions. Taking the one-dimensional case as an example, and writing the position of node i along the line as x_i , the sum of the squared distances spanned by all m edges is

$$\Delta^2 = \frac{1}{2} \sum_{ij} A_{ij}(x_i - x_j)^2, \quad (14.71)$$

where the factor of $\frac{1}{2}$ compensates for the double counting of node pairs.

Upon reflection, however, it is clear that minimizing this quantity alone will not do us any good. Its minimum value of zero is clearly achieved when we set

all x_i equal to the same value, which does not give us a useful embedding. In hindsight this is obvious: the shortest average distance between node pairs is achieved when all of the nodes are in the same place. Moreover, Eq. (14.71) has another problem too: it does not fix the overall position of the nodes in space. If we add a constant to all of the x_i the value of Δ^2 is unchanged, so minimizing it is never going to give us a unique solution for x_i .

We can cure these problems by applying additional constraints to the nodes to fix their overall position in space and to make sure they stay spread out and do not all end up in the same place. There are various ways to do this, but a common choice to fix the overall position is to require that the center of mass of the nodes be at the origin:

$$\sum_i x_i = 0, \quad (14.72)$$

and we can prevent the nodes from all ending up in the same position by fixing the sum of the squares of the positions to have any non-zero value, such as 1 for instance:

$$\sum_i x_i^2 = 1. \quad (14.73)$$

This is equivalent to fixing the variance of the positions, since the mean is zero by Eq. (14.72).

Our goal now is to find the set of positions x_i that minimize Δ^2 , while respecting the constraints (14.72) and (14.73). It is useful first to simplify Eq. (14.71) a little, thus:

$$\begin{aligned} \Delta^2 &= \frac{1}{2} \sum_{ij} A_{ij} (x_i - x_j)^2 = \frac{1}{2} \sum_{ij} A_{ij} (x_i^2 - 2x_i x_j + x_j^2) \\ &= \frac{1}{2} \sum_i k_i x_i^2 - \sum_{ij} A_{ij} x_i x_j + \frac{1}{2} \sum_j k_j x_j^2 = \sum_{ij} (k_i \delta_{ij} - A_{ij}) x_i x_j \\ &= \sum_{ij} L_{ij} x_i x_j, \end{aligned} \quad (14.74)$$

where we have made use of the fact that $\sum_j A_{ij} = k_i$ and $L_{ij} = k_i \delta_{ij} - A_{ij}$ is an element of the graph Laplacian matrix (see Section 6.14). Now we minimize Δ^2 by differentiating with respect to each variable x_v in turn, while imposing the constraints (14.72) and (14.73) using Lagrange multipliers μ and λ :

$$\frac{\partial}{\partial x_v} \left[\sum_{ij} L_{ij} x_i x_j + \mu \sum_i x_i + \lambda \left(1 - \sum_i x_i^2 \right) \right] = 0. \quad (14.75)$$

Performing the derivative, we then find that

$$\sum_j L_{vj}x_j + \frac{1}{2}\mu - \lambda x_v = 0. \quad (14.76)$$

Summing this expression over v gives

$$\sum_{vj} L_{vj}x_j + \frac{1}{2}n\mu - \lambda \sum_v x_v = 0, \quad (14.77)$$

but $\sum_v x_v = 0$ by Eq. (14.72) and

$$\sum_v L_{vj} = \sum_v A_{vj} - \sum_v k_v \delta_{vj} = k_j - k_j = 0, \quad (14.78)$$

and hence (14.77) reduces to $\mu = 0$. Then Eq. (14.76) simplifies to $\sum_j L_{vj}x_j - \lambda x_v = 0$ or, in vector notation,

$$\mathbf{Lx} = \lambda \mathbf{x}, \quad (14.79)$$

where \mathbf{x} is the vector with elements x_i . In other words, the positions x_i of the nodes in the optimal embedding of the network are the elements of an eigenvector of the graph Laplacian, with λ being the corresponding eigenvalue.

Which eigenvector should we use? Substituting $\sum_j L_{ij}x_j = \lambda x_i$ into Eq. (14.74), we get an expression for the optimal value of Δ^2 thus:

$$\Delta^2 = \lambda \sum_i x_i^2 = \lambda, \quad (14.80)$$

where we have used Eq. (14.73). In order to minimize Δ^2 , therefore, we should choose the eigenvector \mathbf{x} corresponding to the smallest eigenvalue λ .

As shown in Section 6.14.5, the smallest eigenvalue of the graph Laplacian is always zero, but the corresponding eigenvector is the uniform vector $\mathbf{x} = (1, 1, 1, \dots)$, which fails to satisfy the constraint (14.72). Thus the best we can do is to choose the eigenvector corresponding to the second-smallest eigenvalue of the Laplacian. The elements of this vector, sometimes called the *Fiedler vector*,¹⁸ tell us the optimal embedding of the network.

As an example, Fig. 14.17 shows the network of interstate highways in the continental United States. The nodes are positioned as they actually fall on the map, and the colors of the nodes reflect the values of the corresponding elements of the Fiedler vector. As we can see, the eigenvector elements vary

¹⁸After the Czech mathematician Miroslav Fiedler, who was one of the first to study its properties.



Figure 14.17: Fiedler vector for the network of US interstate highways. The edges in this network represent the interstate highways of the continental United States and the nodes represent their intersections. The shades of the nodes indicate the values of the elements of the Fiedler vector—the eigenvector of the graph Laplacian corresponding to the second smallest eigenvalue. As we can see, the vector elements vary smoothly from one side of the country to the other, showing that the eigenvector can accurately pick out the spatial embedding of the network.

smoothly from the left of the map to the right, indicating that the eigenvector has successfully identified the true spatial embedding of the network (or at least its east-west dimension) based only on the topology of the network. If we did not already know the spatial structure of the interstate network, we could guess a large part of it from this calculation.

One can extend these methods to embeddings in higher dimensions. Following similar lines of argument one can show that the optimal embedding in a space of d dimensions is to take the d eigenvectors corresponding to the d lowest non-zero eigenvalues. The first elements of each vector, taken together, give the coordinates in the latent space of the first node, the second elements give the coordinates of the second node, and so forth.

The developments here are reminiscent of those of Section 6.14.2 on network visualization, and this is no coincidence. One way to think about network visualization is as an embedding problem. A visualization of a network is precisely a positioning of the nodes in a low-dimensional space (usually the

two-dimensional screen of a computer or a piece of paper), and one way to make a good visualization is to minimize the average lengths of the edges, so that nodes that are network neighbors also tend to be geographic neighbors in the space of the visualization.

Other methods have also been proposed for detecting latent-space structure, although none is as often used as the Laplacian method above. One can address the problem using statistical inference, for instance [235,307,368]. One defines a model of a network in a latent space, then fits it to an observed network using maximum likelihood methods similar to those of Section 14.4.1. Consider again the example of a one-dimensional latent space. In the typical model we would assign to each node a position in the space, then place edges between node pairs with some probability $\omega(x, y)$, which is a function of the positions x, y of the nodes. The model can show different behaviors depending on the form of the function $\omega(x, y)$. Hoff *et al.* [235], for instance, assume a function of the form

$$\omega(x, y) = \frac{ce^{-|x-y|}}{1 + ce^{-|x-y|}} \quad (14.81)$$

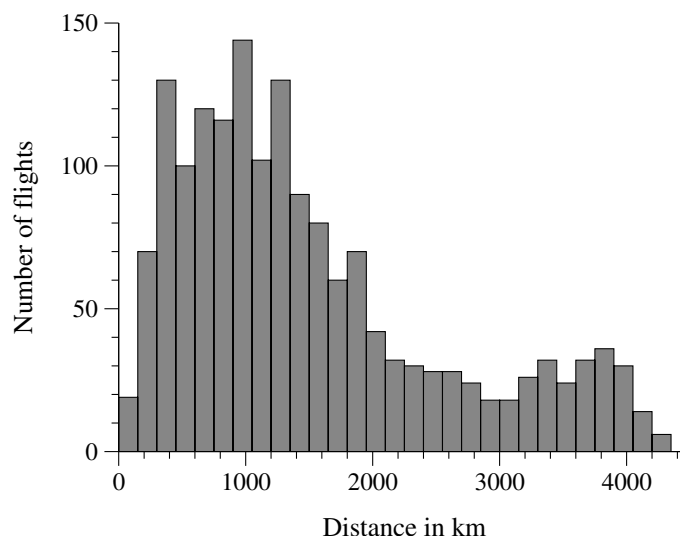
for some value of the constant c . This function takes the value $c/(1 + c)$ when x and y are in the same position and decreases monotonically as x and y move apart. Thus this form assumes a larger probability of connection between nodes that are close together than nodes far apart.

One of the nice features of the inference method, however, is that one does not need to make such an assumption if one does not want to. What if one had a network where nodes were more likely to be connected if they were far apart? Or there could be some more complicated dependence of edge probability on distance. Consider Fig. 14.18, for instance, which shows the distribution of lengths of edges in a network of US airline flights. As we can see, the distribution is non-monotone, with very little weight at the shortest distances (because airlines don't fly between very closely spaced airports), rising to a peak around 1000 km, falling to a minimum around 3000 km, and then rising to another peak around 4000 km (representing coast-to-coast flights). Quite general forms of the probability function $\omega(x, y)$, capable of capturing non-trivial connection patterns like this one, can be specified, either by expressing the function as a linear combination of a suitable set of basis states, or by using “non-parametric” methods that allow any functional form subject to basic conditions of smoothness. Then the process of fitting the model becomes one of choosing the best form from the proposed family of possibilities—see Refs. [307,368] for examples.

Latent-space structure can also occur in directed networks, where both the existence of edges and their direction can depend on node positions in the latent

In statistics this is called the *logistic function* or *inverse-logit*. In physics it is called the *Fermi function*.

Figure 14.18: Distribution of lengths of US airline flights. A histogram of the (straight line) distances covered by passenger flights operated by one major airline in the continental United States. The distribution is clearly non-monotonic, having two separate peaks, one around 1000 km, which accounts for the bulk of the flights, and a second smaller one around 4000 km, which corresponds to coast-to-coast flights. After Gastner and Newman [203].



This use of the word hierarchy is different in meaning from that of either Section 14.5.2 or Section 14.7.2, and one should be careful to keep the distinction clear.

See Section 6.4.1 for an introduction to acyclic networks.

space. There has been relatively little work on directed networks, however, except for one special case, the case of an *ordered network*, also called *hierarchy* or *rank structure*. Rank structure occurs when nodes lie in a one-dimensional latent space and all or most of the directed edges between them point in the same direction in that space, resulting in a network that is acyclic or mostly acyclic. We have seen a number of examples in our previous discussions of acyclic networks. Citation networks, for instance, can be thought of as rank-structured networks where the latent space is time, and food webs can be thought of as networks where the latent space is trophic level. In a citation network the latent-space positions of the nodes (i.e., their publication dates) are usually known, but for other networks, such as food webs, they may not be, in which case we can use methods like statistical inference to determine them [38]. Another approach is the so-called *minimum violations ranking* [19], in which one attempts to find the ranking of the nodes that maximizes the number of edges that point in the same direction (and minimizes the number that point in the opposite direction—hence “minimum violations”).¹⁹

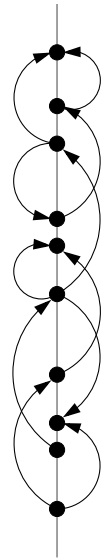
As an example, Ball and Newman [38] used a maximum likelihood method to infer rank structure in networks of directed friendships among high school students, finding that there is a clear ranking of students from high to low

¹⁹Centrality measures such as PageRank (Section 7.1.4) are also aimed at ranking the nodes of directed networks, although they do not explicitly try to make all edges point in the same direction. See Ref. [205] for a discussion of the use of PageRank for general network ranking problems.

such that lower ranked students claim to be friends with higher ranked ones, but the reverse rarely happens. Ball and Newman speculated that this ranking corresponds, at least roughly, to some type of social status or pecking order within schools. In another study, Clauset *et al.* [107] calculated a minimum violations ranking for a network of faculty hiring at US universities. The nodes in this network represent universities and the directed edges represent individuals who received a PhD at one university and then found a job on the faculty of another. Clauset *et al.* found a clear ranking of universities from high to low such that most edges pointed downward—most people found employment at a university ranked lower than the one from which they received their PhD. The ranking so derived, moreover, corresponds quite closely to rankings of US universities given by independent observers.

A practical application of these kinds of ranking methods is to the grading of competitions such as sporting events or leagues [289]. One can represent the pattern of games between competitors in a league as a network where the nodes represent the competitors and directed edges represent games or matches between them, with the direction of each edge pointing from the winner to the loser of the corresponding game. If there is an unambiguous ranking of competitors from best to worst, and if the better of two competitors always wins when they play a game, then the resulting network will be perfectly acyclic: if the competitors are arranged in a line in rank order then all edges will point from higher to lower ranked nodes. But if we don't know the ranking of the competitors then we can use the network to estimate it by the kind of methods described earlier. For instance, we can apply an inference method similar to that of Ref. [38] or find a minimum violations ranking that ranks as many competitors as possible higher than the opponents they defeated [380].

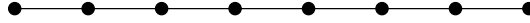
In another context, such methods are also used to rank animals in the networks of aggression known as dominance hierarchies (see Section 4.3). In these networks, animals engage in fights or other aggressive behaviors in an attempt to establish social dominance over one another [150], and the resulting patterns of wins and losses are in many ways analogous to those of sports tournaments and can be analyzed in essentially the same way [136].



A network with rank structure, in which nodes are positioned in a one-dimensional latent space represented by the vertical line, and most (though in this case not all) edges point from lower to higher nodes.

EXERCISES

14.1 Consider a “line graph” consisting of n nodes in a row like this:

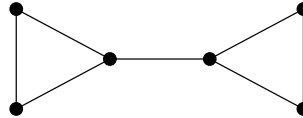


a) Show that if we divide the network into two parts by cutting any single edge, such that one part has r nodes and the other has $n - r$, the modularity, Eq. (7.58), takes the value

$$Q = \frac{3 - 4n + 4rn - 4r^2}{2(n - 1)^2}.$$

b) Hence show that when n is even the optimal such division, in terms of modularity, is the division that splits the network exactly down the middle.

14.2 Using your favorite numerical software for matrices construct the modularity matrix for this small network:



Find the eigenvector of the modularity matrix corresponding to the largest eigenvalue and hence divide the network into two communities.

14.3 Suppose we have a set of n integers k_1, \dots, k_n that are drawn independently at random from a Poisson distribution with mean μ . In other words, the likelihood of drawing an integer k is

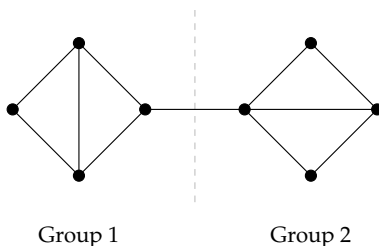
$$P(k|\mu) = \frac{\mu^k}{k!} e^{-\mu}. \tag{14.82}$$

- a) Given that our n numbers are statistically independent, derive an expression for the log-likelihood $\mathcal{L} = \log P(k_1, \dots, k_n|\mu)$ of the complete set of values k_1, \dots, k_n .
- b) Suppose we are given only the observed values k_i and we are told that they are drawn from a Poisson distribution, but we are not told the value of μ . Derive an expression for the best estimate of μ by maximizing the log-likelihood.

14.4 Consider a random graph model of a bipartite network in which there are n_1 nodes of one type, n_2 nodes of another type, and every pair of nodes of unlike types is connected by an edge independently with probability p or not with probability $1 - p$.

- a) Let \mathbf{B} be the $n_1 \times n_2$ incidence matrix of an observed bipartite network. Derive an expression for the likelihood (i.e., the probability) that this network was generated by the model for a given value of p . Take the log of your expression to find the log-likelihood.
- b) Find the value of p that maximizes the log-likelihood.

14.5 Consider this small network, divided into two groups as indicated:



- For this network calculate the three quantities m_{rs} and the two quantities κ_r that appear in the (log) profile likelihood, Eq. (14.50), for the degree-corrected stochastic block model. Hence calculate the numerical value of the profile likelihood.
- Verify that no higher profile likelihood can be achieved by moving any single node to the other group, and hence that this division into groups is at least a local maximum of the likelihood. (In fact it's the global maximum as well.) Hint: Some of the nodes are symmetry equivalent, which means you need only consider the movement of six different nodes to the other group, which could save you some effort.

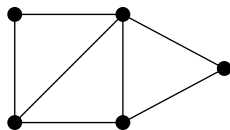
14.6 Consider the following network model. Each of n nodes is assigned a non-negative real parameter θ_i and then undirected edges are created such that the number of edges between nodes i and j is a Poisson-distributed independent random number with mean $\theta_i\theta_j$, except for self-edges, which have mean $\frac{1}{2}\theta_i^2$. The goal is to find the values of the parameters θ_i that best fit a given observed network.

- Derive an expression for the likelihood (i.e., the probability) that a network with adjacency matrix \mathbf{A} was generated by this model, given the values of the parameters θ_i , and hence show that the log-likelihood $\mathcal{L} = \log P(\mathbf{A}|\boldsymbol{\theta})$, ignoring constants that don't depend on θ_i , is

$$\mathcal{L} = \frac{1}{2} \sum_{ij} [A_{ij} \log(\theta_i\theta_j) - \theta_i\theta_j].$$

- By maximizing the log-likelihood with respect to the θ parameters show that for the best fit of this model to an observed network, the mean number of edges between distinct nodes i and j is $k_i k_j / 2m$, where k_i is the observed degree of node i and m is the number of edges in the network. (In other words, this model is basically the configuration model.)

14.7 Consider this small network with five nodes:



COMMUNITY STRUCTURE

- a) Calculate the cosine similarity for each of the $\binom{5}{2} = 10$ pairs of nodes. (For a definition of cosine similarity see Section 7.6.1.)
- b) Using the values of the ten similarities construct the dendrogram for the single-linkage hierarchical clustering of the network according to cosine similarity.

CHAPTER 15

PERCOLATION AND NETWORK RESILIENCE

A discussion of one of the simplest of processes taking place on networks, percolation, and its use as a model of network resilience

STUDIES of the structure of networks, which have been our primary focus in this book so far, are only one step towards the understanding of networked systems. Another important part of the puzzle is to make the connection between network structure and function: once we have measured and quantified the structure of a network, how do we turn our new knowledge into predictions or conclusions about how the overall system will behave? Unfortunately, we do not yet have a comprehensive theory of the connection between structure and function in networks, but there are a number of individual areas in which progress has been made, several of which we examine in the next few chapters. In this chapter we study the phenomenon of percolation, which leads to an elegant theory of the robustness of networked systems to the failure of their components.

15.1 PERCOLATION

Imagine taking a network and removing some fraction of its nodes, along with the edges connected to those nodes—see Fig. 15.1. This process is called *percolation* and it can be used as a model of a variety of real-world phenomena. The failure of routers on the Internet, for instance, can be formally represented by removing the corresponding nodes and their attached edges from a network

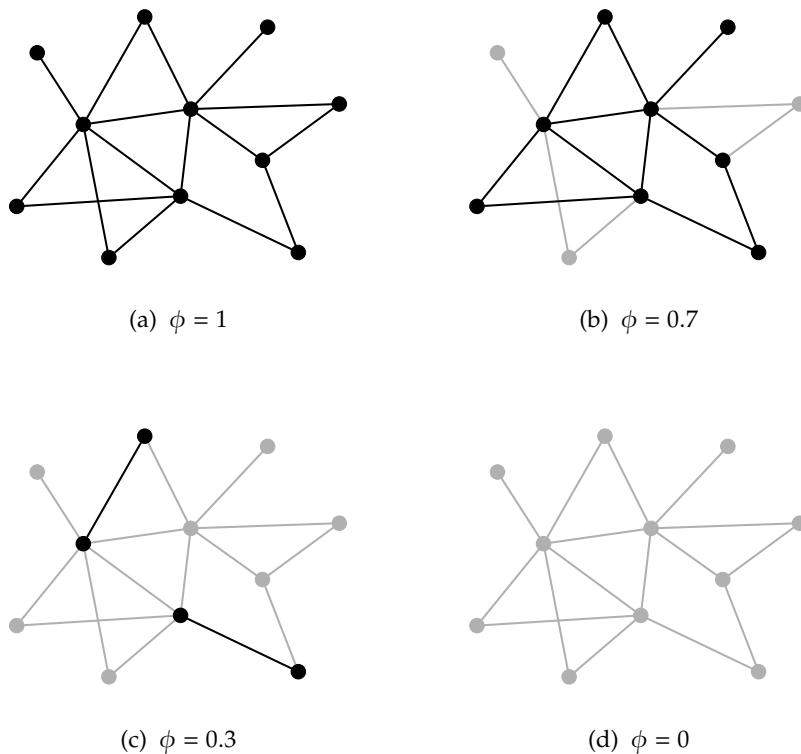


Figure 15.1: Percolation. A depiction of the site percolation process on a small network for various values of the occupation probability ϕ . Gray denotes nodes that have been removed, along with their associated edges, and black denotes those that are still present. The networks in (a) and (b) are above the percolation threshold while those in (c) and (d) are below it.

representation of the Internet. In fact, about 3% of the routers on the Internet are non-functional for one reason or another at any one time, and it is a question of some practical interest what effect this will have on the performance of the network. The theory of percolation processes can help us answer this question.

Another example of a percolation process is the vaccination or immunization of individuals against the spread of a disease. As discussed in Chapter 1, and at greater length in Chapter 16, diseases spread through populations over the networks of contacts between individuals. But if an individual is vaccinated against a disease and therefore cannot catch it, then that individual does not

contribute to the spread of the disease. The individual is still present in the network, but, from the point of view of the spread of the disease, might as well be absent, and hence the vaccination process can again be formally represented by removing nodes.

One can see immediately that percolation processes can give rise to some interesting behaviors. The vaccination of an individual in a population, for example, not only prevents that individual from becoming infected but also prevents them from infecting others, and so has a “knock-on” effect in which the benefit of vaccinating one individual is felt by more than one. As we will show, this knock-on effect means that in some cases the vaccination of a relatively small fraction of the population can effectively prevent the spread of disease to anyone, an outcome known as *herd immunity*.

Similar effects crop up in our Internet example, although in that case they are usually undesirable. The removal or failure of a single router on the Internet prevents that router from receiving data, but also prevents data from reaching others via the failed router, forcing traffic to take another path—possibly longer or more congested—or even cutting off some portions of the network altogether. One of the goals of percolation theory on networks is to understand how the knock-on effects of node removal or failure affect the network as a whole.

Sometimes it is not the nodes in the network that fail but the edges. For instance, communication lines on the Internet can fail, disconnecting routers from one another, even though the routers themselves may still be functioning perfectly. Phenomena like this can be modeled using a slightly different percolation process in which edges rather than nodes are removed from the appropriate network. If we need to distinguish between the two types of percolation process we could refer to them as node percolation on the one hand and edge percolation on the other, but in fact they are more commonly called *site percolation* and *bond percolation*, a nomenclature that derives from physics.¹ In this chapter we will focus principally on site percolation (i.e., removal of nodes) but bond percolation (removal of edges) will become important in Chapter 16 when we look at epidemic processes.

15.2 UNIFORM RANDOM REMOVAL OF NODES

There is more than one way in which nodes can be removed from a network. In the simplest case they could be removed purely at random: we could, for

¹If you are interested in the study of percolation in physics, the book by Stauffer and Aharony [437] contains a lot of interesting material on the subject, although most of it is not directly relevant to networks.

example, take away some specified fraction of the nodes chosen uniformly at random from the entire network. This is the most commonly studied form of site percolation, and indeed for many people the word “percolation” refers specifically to this process. But there are other ways in which nodes could be removed, and percolation as studied in this chapter is considered to include all of them. One popular alternative removal scheme is to remove nodes according to their degree in some fashion. For instance, we could remove nodes in order of degree from highest to lowest, an approach that turns out to make an effective vaccination strategy for the control of disease. Other approaches have also been considered occasionally, such as removing nodes with high betweenness centrality. Let us begin, however, by examining the simple case of uniformly random removal.

Consider a network in which some fraction of the nodes, selected uniformly at random, are removed. As discussed in Section 15.1, in many real-world situations “removal” does not imply actual physical removal of the nodes, only that they are non-functional in some way, such as routers that have failed on the Internet, or vaccinated individuals in a network of disease-causing contacts.

Traditionally the percolation process is parametrized by a quantity ϕ equal to the probability that a node is present or functioning in the network. In the parlance of percolation theory, one says that the functional nodes are *occupied* and ϕ is called the *occupation probability*. Thus $\phi = 1$ indicates that all nodes in the network are occupied (i.e., no nodes have been removed) and $\phi = 0$ indicates that no nodes are occupied (i.e., all of them have been removed).²

Now look again at Fig. 15.1 and consider panel (a), in which $\phi = 1$, all nodes are present or occupied, and all nodes are connected together into a single component. (The network could have more than one component, but in this example it has only one.) Now look at the other panels. In panel (b) a few nodes have been removed, but those that remain are all still connected together by the remaining edges. In panel (c) still more nodes have been removed, and now so many are gone that the few remaining nodes are no longer all connected together, but divide into two small components. In the final panel, panel (d), all nodes have been removed and there is no network left at all.

The behavior of this small example is typical of percolation processes. When ϕ is large the nodes tend to be connected together, forming a giant component that fills most of the network (although there may be small components also). But as ϕ is decreased there comes a point where the giant component breaks

²In most of the physics literature on percolation the occupation probability is denoted p . We, however, will use ϕ because the letter p is used for many other things in the theory of networks and could cause confusion.

apart and we are left only with small components. Conversely, if we increase ϕ from zero we first form small components, which grow in size and eventually coalesce to form a giant component that fills a large fraction of the network.

The formation or dissolution of a giant component in this fashion is called a *percolation transition*. When the network contains a giant component we say that it *percolates* and the point at which the percolation transition occurs is called the *percolation threshold*.

The percolation transition is similar in many ways to the phase transition in the Poisson random graph at which a giant component forms (see Section 11.5). In the random graph we vary not the fraction of occupied nodes but the probability of connection between those nodes. In both cases, however, when enough of the network is removed, the giant component is destroyed and we are left with only small components.

In studies of percolation, the “components” that remain after nodes have been removed are in fact usually called *clusters*, another term inherited from the physics literature and one that we will use here—it will be useful to distinguish between the “components” of the underlying network and the “clusters” of the percolation process. That is, we will use “component” to refer to connected groups of nodes on the original network before any nodes have been removed and “cluster” to refer to those after removal. The giant component of the percolation process, if there is one, is thus properly called the *giant cluster* or sometimes the *percolating cluster*.³

The percolation transition plays a central role in our interpretation of percolation phenomena. In a network like the Internet, for example, there has to be a giant cluster if the network is to perform its intended function as a communications network. If the network has only small clusters, as in Fig. 15.1c, then every node has a connection to, at most, a handful of others and is cut off from everyone else. If there is a giant cluster, on the other hand, then the members of that giant cluster, who are a non-zero fraction of all nodes in the network, are connected and can communicate with one another, although the remainder of the network is still cut off. Thus, the presence of a giant cluster is an indicator of a network that is at least partly performing its intended function, while the size of the giant cluster tells us exactly how much of the network is working.

³Elsewhere in the literature on percolation the giant cluster is also called the *spanning cluster*. The reason for this name is that most work on percolation has been on low-dimensional lattices such as the square lattice. On such lattices the giant cluster is distinguished by being the only cluster that spans the lattice from one side to the other in the limit of large n . There is no equivalent phenomenon for percolation on general networks, however, since networks don’t have sides, so the concept of spanning is not a useful one.

15.2.1 UNIFORM REMOVAL IN THE CONFIGURATION MODEL

To gain some understanding of the percolation transition and the giant cluster, let us consider the behavior of the site percolation process on networks generated using the configuration model of Chapter 12, a simple but useful model of a network with a specified degree distribution. We can calculate the properties of the giant percolation cluster in the configuration model by a method similar to the one we used for the giant component in Section 12.6.

Consider a configuration model network with degree distribution p_k and a percolation process on that network in which nodes are present or occupied with occupation probability ϕ as described in Section 15.2. Now consider one of the nodes that is present in the network (i.e., one that has not been removed). If that node is to belong to the giant cluster it must be connected to it via at least one of its neighbors. Equivalently, it is not a member of the giant cluster if and only if it is not connected to the giant cluster via any of its neighbors. Let us define u to be the average probability that a node is not connected to the giant cluster via a particular neighbor. Then if the node in question has degree k , the total probability of its not belonging to the giant cluster is u^k . And if we then average over the probability distribution p_k of the degree we find that the average probability of not being in the giant cluster is $\sum_k p_k u^k = g_0(u)$, where

$$g_0(z) = \sum_{k=0}^{\infty} p_k z^k \quad (15.1)$$

is the generating function for the degree distribution, as defined previously in Eq. (12.96). Then the average probability that a node *does* belong to the giant cluster is $1 - g_0(u)$.

Bear in mind, however, that this is for a node that is itself assumed not to have been removed from the network. Nodes that have been removed are obviously not members of the giant cluster. Thus, out of all the original nodes in the network, the total fraction S that are in the giant cluster is equal to the fraction ϕ that have not been removed times the probability $1 - g_0(u)$ that they are in the giant cluster:

$$S = \phi[1 - g_0(u)]. \quad (15.2)$$

To use this equation we still need to calculate the value of u , which is the average probability that a node is not connected to the giant cluster via a particular neighboring node. There are two ways to not be connected to the giant cluster via a neighbor: either the neighbor in question—let us call it node i —has been removed, which happens with probability $1 - \phi$, or it is present (probability ϕ) but it is not itself a member of the giant cluster. The

latter happens if and only if i is not connected to the giant cluster via any of its other neighbors, which occurs with probability u^k where k is the number of other neighbors. Adding everything together, the total probability that we are not connected to the giant cluster via i is $1 - \phi + \phi u^k$.

Since i is reached by following an edge, k in this case is the excess degree of i and obeys the excess degree distribution

$$q_k = \frac{(k+1)p_{k+1}}{\langle k \rangle}, \tag{15.3}$$

See Section 12.2 for a discussion of the excess degree distribution.

where $\langle k \rangle$ is the average degree in the network. Averaging over this distribution, we then arrive at an expression for the average probability u thus:

$$\begin{aligned} u &= \sum_{k=0}^{\infty} q_k (1 - \phi + \phi u^k) = 1 - \phi + \phi \sum_{k=0}^{\infty} q_k u^k \\ &= 1 - \phi + \phi g_1(u), \end{aligned} \tag{15.4}$$

where

$$g_1(z) = \sum_{k=0}^{\infty} q_k z^k \tag{15.5}$$

is the generating function for the excess degree distribution, defined previously in Eq. (12.97), and we have made use of the normalization condition $\sum_k q_k = 1$. Together, Eqs. (15.2) and (15.4) give us a complete solution for the size of the giant cluster in our network.⁴

In practice it is often not possible to solve Eq. (15.4) in closed form, but there is an elegant graphical representation of the solution as follows. Consider Fig. 15.2a, which gives a sketch of the function $g_1(u)$. The exact form of the curve will depend on the degree distribution, but we know the general shape: g_1 is a polynomial with all coefficients non-negative (because they are probabilities), so for $u \geq 0$ it must have a non-negative value and all derivatives non-negative. Thus in general it is an increasing function of u and curves upward as shown in the figure.

To get the function $1 - \phi + \phi g_1(u)$ that appears on the right-hand side of Eq. (15.4) we first multiply $g_1(u)$ by ϕ then add $1 - \phi$. Graphically, this is equivalent to compressing the unit square of Fig. 15.2a, along with the curve it

⁴This solution of the percolation problem has a history stretching back some decades. In 1961, Fisher and Essam [179] derived a solution for percolation on regular trees (called Cayley trees or Bethe lattices in physics), which is equivalent to the solution given here for the case where every node has the same degree. The developments for general degree distributions, however, were not given till much later [93, 113].

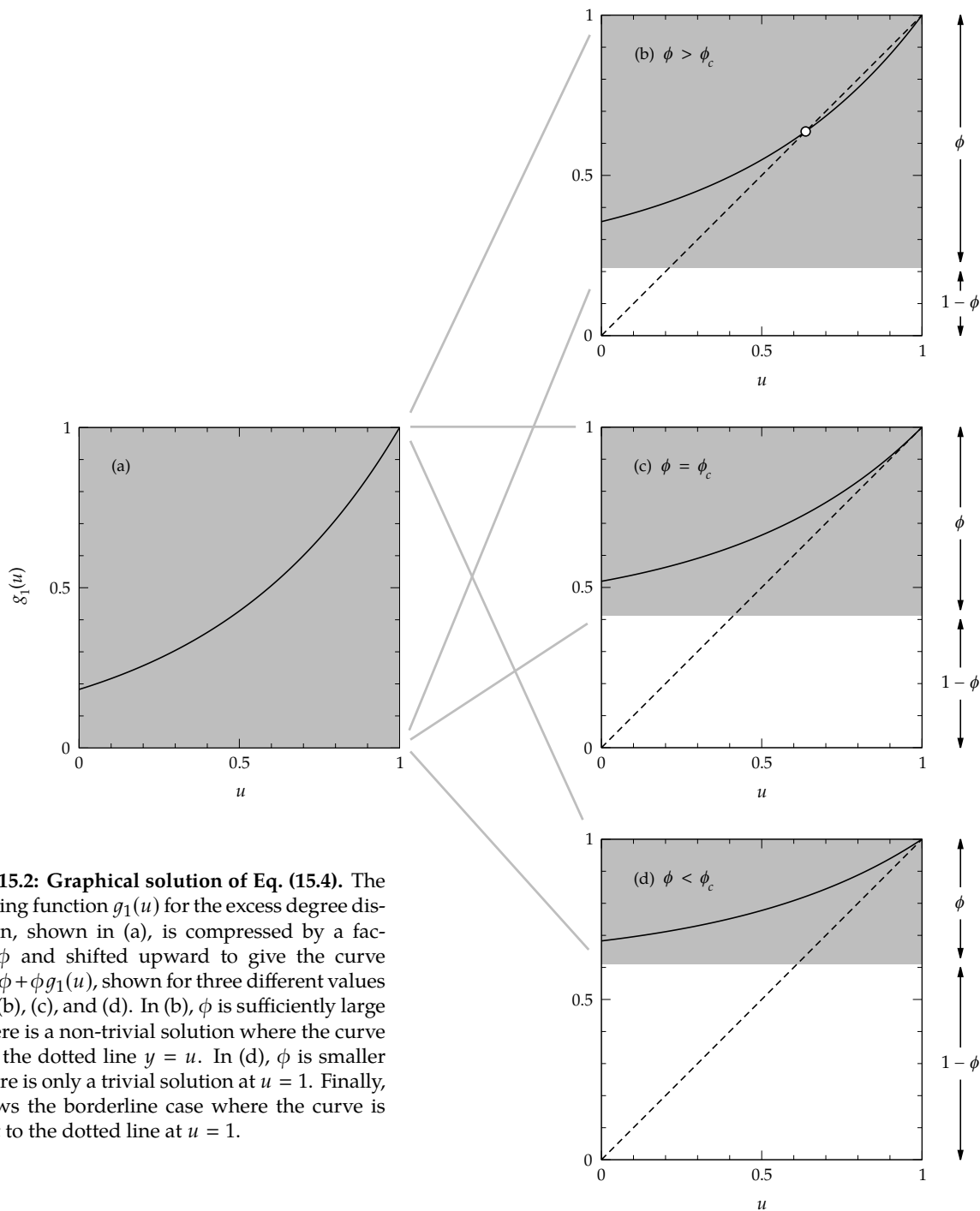


Figure 15.2: Graphical solution of Eq. (15.4). The generating function $g_1(u)$ for the excess degree distribution, shown in (a), is compressed by a factor of ϕ and shifted upward to give the curve $y = 1 - \phi + \phi g_1(u)$, shown for three different values of ϕ in (b), (c), and (d). In (b), ϕ is sufficiently large that there is a non-trivial solution where the curve crosses the dotted line $y = u$. In (d), ϕ is smaller and there is only a trivial solution at $u = 1$. Finally, (c) shows the borderline case where the curve is tangent to the dotted line at $u = 1$.

contains, until it has height ϕ and then shifting it upward a distance $1 - \phi$ as shown in Fig. 15.2b. The point or points at which the resulting curve crosses the line $y = u$ (dotted line in Fig. 15.2b) are then the solutions to Eq. (15.4).

In Fig. 15.2b there are two such solutions. One is a trivial solution at $u = 1$. This solution always exists because $g_1(1) = 1$ for any correctly normalized excess degree distribution q_k . But there is also a non-trivial solution with $u < 1$, indicated by the dot in the figure. Only if we have such a non-trivial solution can there be a giant cluster in the network and the value of u for this solution gives us the size of the giant cluster via Eq. (15.2). (The $u = 1$ solution gives $S = 0$ in Eq. (15.2) and so doesn't give us a giant cluster.)

Now consider Fig. 15.2d, which shows the solution of Eq. (15.4) for a smaller value of ϕ . Now the curve of the generating function has been compressed more and the result is that the non-trivial solution for u has vanished. Only the trivial solution at $u = 1$ remains and so in this regime there can be no giant cluster.

Figure 15.2c shows the borderline case between cases (b) and (d). The non-trivial solution for u vanishes at this point, where the curve just meets the dotted line. This is the percolation threshold. Mathematically, this is the point at which the curve is tangent to the dotted line at $u = 1$, i.e., the point where its gradient at $u = 1$ is 1. In other words, the percolation threshold occurs when

$$\left[\frac{d}{du}(1 - \phi + \phi g_1(u)) \right]_{u=1} = 1. \quad (15.6)$$

Performing the derivative and rearranging, we then find that the value of ϕ at the transition, which we call the *critical value* and denote ϕ_c , is

$$\phi_c = \frac{1}{g_1'(1)}. \quad (15.7)$$

We can express the critical value more directly in terms of the degree distribution by making use of the definition of the generating function $g_1(z)$ and the excess degree distribution in Eqs. (15.3) and (15.5). Substituting one into the other and differentiating, we find that

$$\begin{aligned} g_1'(1) &= \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} k(k+1)p_{k+1} = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} k(k-1)p_k \\ &= \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k \rangle}, \end{aligned} \quad (15.8)$$

and hence the critical occupation probability ϕ_c is given by

$$\phi_c = \frac{\langle k \rangle}{\langle k^2 \rangle - \langle k \rangle}, \quad (15.9)$$

an expression first given by Cohen *et al.* [113].

This equation tells us the minimum fraction of nodes that must be present or occupied in our configuration model network for a giant cluster to exist. Thus, for instance, if we were to consider the configuration model as a simple model of the Internet, we would want to make ϕ_c low, so that the network will have a giant cluster even when a significant fraction of nodes are non-functional, and hence go on functioning as a communication network. If, for example, the network had a Poisson degree distribution,

$$p_k = e^{-c} \frac{c^k}{k!}, \quad (15.10)$$

where c is the mean degree, then $\langle k \rangle = c$ and $\langle k^2 \rangle = c(c + 1)$, so

$$\phi_c = \frac{1}{c}. \quad (15.11)$$

Then if we make c large we will have a network that can withstand the loss of many of its nodes. For $c = 4$, for example, we would have $\phi_c = \frac{1}{4}$, meaning that $\frac{3}{4}$ of the nodes would have to fail before the giant cluster is destroyed. A network that can tolerate the loss of a large fraction of its nodes in this way is said to be *robust* against random failure.

The degree distribution of the Internet, however, is not Poissonian. In fact, as discussed in Section 10.4, the Internet's degree distribution appears roughly to follow a power law with an exponent $\alpha \simeq 2.5$ (see Table 10.1). As we showed in Section 10.4.2, power laws with exponents in the range $2 < \alpha < 3$, which includes most real-world examples, have a finite mean $\langle k \rangle$, but their second moment $\langle k^2 \rangle$ diverges. In this case Eq. (15.9) implies that $\phi_c = 0$. In other words, no matter how many nodes we remove from the network there will *always* be a giant cluster. Scale-free networks—those with power-law degree distributions—are thus highly robust networks that can survive the failure of any number of their nodes, a point first highlighted in the work of Albert *et al.* [17].

In practice, as discussed in Section 10.4.2, the second moment of the degree distribution is never actually infinite in any finite network. It can still become very large, however, which can result in non-zero but very small values of ϕ_c , so that the network is still highly robust.

The structure of the real Internet is not the same as that of a configuration model. It has many layers and levels of structure engineered into it, as discussed in Section 2.1. Nonetheless, it does appear to be quite robust to random removal of its nodes. For instance, Albert *et al.* [17] simulated the behavior of the Internet as nodes were randomly removed and found that performance is hardly affected

at all by the removal of even a large fraction of the nodes. (Performance is of course completely destroyed for the nodes that are themselves removed, but for the remaining ones the effects are relatively minor.) Thus, as is often the case, our simple network model gives us a good general guide to the behavior of the system, even if it gets some of the details wrong.

Network robustness also plays an important role in the vaccination example mentioned at the start of the chapter. A disease spreading over a contact network between individuals can only reach a significant fraction of the population if there is a giant cluster in the network. If the network contains only small clusters, then an outbreak of the disease will be hemmed in and unable to spread beyond the small cluster in which it starts. Thus one does not have to vaccinate the entire population to prevent disease spread. One need only vaccinate enough of them to bring the network below its percolation threshold and eliminate the giant cluster. This is the herd immunity effect we mentioned earlier.

In the vaccination case, network robustness is a bad thing. The fewer individuals we have to vaccinate to destroy the giant cluster the better. Thus small values of ϕ_c are now bad and large values are good. Unfortunately, we usually don't have much control over the degree distributions of contact networks, so we may be stuck with a low value of ϕ_c whether we like it or not. In particular, if the network in question has a power-law (or approximately power-law) degree distribution, then ϕ_c may be very small, implying that almost all nodes have to be vaccinated to wipe out the disease. Some contact networks do indeed appear to have roughly power-law degree distributions [253, 304, 305] and it may be very difficult to eradicate some diseases as a result [383].

It is interesting to ask how the special behavior of power-law networks shows up in the graphical solution of Fig. 15.2. The answer is that, since $g_1'(1)$ is infinite in the power-law case (because $\langle k^2 \rangle$ diverges in Eq. (15.8) while $\langle k \rangle$ remains finite), the curve of $g_1(u)$ has infinite slope at $u = 1$. Thus $g_1(u)$ must look something like Fig. 15.3. Because of the infinite slope, it makes no difference how much we compress the function (as in Fig. 15.2)—the curve will

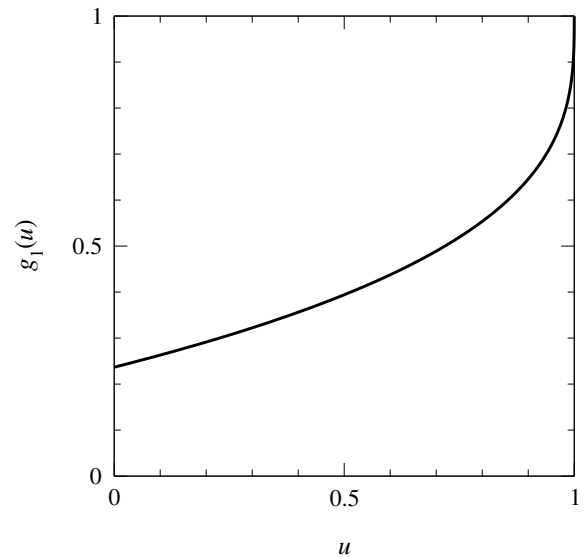


Figure 15.3: Generating function for the excess degree distribution in a scale-free network. The generating function $g_1(u)$ for a network with a power-law degree distribution has a derivative that diverges as $u \rightarrow 1$, though the value of the generating function itself remains finite and tends to 1 in this limit. Thus the function looks generically like the curve sketched here.

always drop below the line of $y = u$ before coming back up again and crossing it, giving a non-trivial solution for u .

The position of the percolation threshold is not the only quantity important in assessing the robustness of a network. The size of the giant cluster also plays a role because it tells us what fraction of the network will be connected and functional. To find the size of the giant cluster we need to solve Eq. (15.4) for u and then substitute the result back into Eq. (15.2). In many cases, as we have said, we cannot solve for u exactly, but in some cases we can. Consider, for example, a network with an exponential (or geometric) degree distribution given by

$$p_k = (1 - a) a^k, \quad (15.12)$$

where $a < 1$ and the leading factor of $1 - a$ ensures that the distribution is properly normalized. Then, as discussed in Section 12.10.6, we have

$$g_0(z) = \frac{1 - a}{1 - az}, \quad g_1(z) = \left(\frac{1 - a}{1 - az} \right)^2, \quad (15.13)$$

and Eq. (15.4) becomes

$$u(1 - au)^2 - (1 - \phi)(1 - au)^2 - \phi(1 - a)^2 = 0. \quad (15.14)$$

This is a cubic equation in u , which is ugly to solve, though not impossible. In this case, however, we don't have to solve it directly. We observe instead that $u = 1$ is always a solution of Eq. (15.4) and hence that our cubic equation must contain a factor of $u - 1$. A few moments' work reveals that indeed this is correct. Equation (15.14) factorizes as

$$(u - 1)[a^2 u^2 + a(a\phi - 2)u + \phi - 2a\phi + 1] = 0. \quad (15.15)$$

Thus the two other solutions for u satisfy the quadratic equation

$$a^2 u^2 + a(a\phi - 2)u + \phi - 2a\phi + 1 = 0. \quad (15.16)$$

Of these two solutions one is greater than 1 for $a < 1$ and so cannot be our probability u . The other is

$$u = a^{-1} - \frac{1}{2}\phi - \sqrt{\frac{1}{4}\phi^2 + \phi(a^{-1} - 1)}. \quad (15.17)$$

Now we can plug this value back into Eq. (15.2) to get an expression for the size

of the giant cluster as a fraction of the whole network:

$$\begin{aligned}
 S &= \phi \left[1 - \frac{2(a^{-1} - 1)}{\phi + \sqrt{\phi^2 + 4\phi(a^{-1} - 1)}} \right] \\
 &= \phi \left[1 - 2(a^{-1} - 1) \frac{\phi - \sqrt{\phi^2 + 4\phi(a^{-1} - 1)}}{\phi^2 - (\phi^2 + 4\phi(a^{-1} - 1))} \right] \\
 &= \frac{3}{2}\phi - \sqrt{\frac{1}{4}\phi^2 + \phi(a^{-1} - 1)}. \tag{15.18}
 \end{aligned}$$

Note that the solution for u , Eq. (15.17), can become greater than 1 for sufficiently small ϕ , which is unphysical. In this regime the only acceptable solution is the trivial $u = 1$ solution, which gives $S = 0$ and so there is no giant cluster when this happens. This gives us an alternative way to derive the position of the percolation transition. The transition takes place at the point where Eq. (15.17) equals one, i.e., when

$$a^{-1} - 1 - \frac{1}{2}\phi = \sqrt{\frac{1}{4}\phi^2 + \phi(a^{-1} - 1)}. \tag{15.19}$$

Squaring both sides and rearranging for ϕ we find that the percolation threshold falls at

$$\phi_c = \frac{1 - a}{2a}. \tag{15.20}$$

It is left as an exercise to demonstrate that this is the same as the result we get if we apply the general formula, Eq. (15.7).

Note also that if a is less than $\frac{1}{3}$, then the value of ϕ_c given by Eq. (15.20) is greater than 1. For values of a this small ϕ can thus never be greater than ϕ_c , so there is no percolation transition and the system never percolates. Upon closer inspection, it turns out that $a = \frac{1}{3}$ is precisely the point at which the network itself loses its giant component,⁵ which explains why percolation is not possible beyond this point. For $a < \frac{1}{3}$ the network has no giant component, and hence it is not possible to have a giant cluster even if every node in the network is occupied. (A similar result applies to all networks of course—a giant percolation cluster is never possible in a network without a giant component.)⁶

⁵From Eq. (12.42) we know that a configuration model network has a giant component if and only if $g'_1(1) > 1$, and thus loses its giant component at the point where $g'_1(1) = 1$. Substituting from Eq. (15.13), our network loses its giant component when $2a/(1 - a) = 1$, i.e., when $a = \frac{1}{3}$. See also Problem 12.11 on page 431 for another derivation of this result.

⁶It's straightforward to see this in the case of the configuration model, where $\phi_c = 1/g'_1(1)$ is greater than 1 whenever $g'_1(1) < 1$, i.e., precisely when the network does not have a giant component. Thus, the point at which ϕ_c exceeds 1 and percolation becomes impossible always coincides with the point at which the giant component disappears.

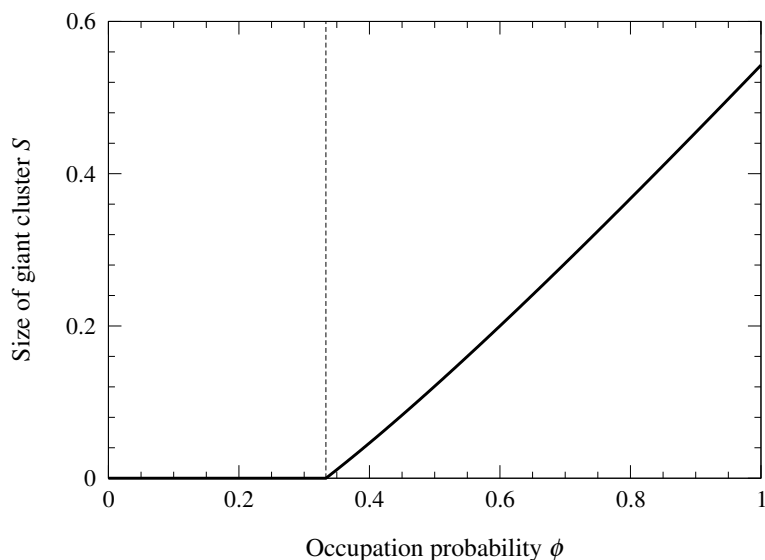


Figure 15.4: Size of the giant cluster for site percolation in the configuration model. The curve indicates the size of the giant cluster for a configuration model with an exponential degree distribution of the form (15.12) with $a = 0.6$, as given by Eq. (15.18). The dashed line indicates the position of the percolation transition, Eq. (15.20), which falls at $\phi_c = \frac{1}{3}$ in this case.

Figure 15.4 shows a plot of the value of S for our exponential network with $a = 0.6$, as a function of ϕ . For small ϕ there is a region in which there are only small clusters and no giant cluster. When we pass through the percolation transition, marked by the dashed line in the figure, a giant cluster appears and grows smoothly from zero as ϕ increases. This is an example of what a physicist would call a continuous phase transition.⁷ We saw other examples in Sections 11.5 (for the Poisson random graph) and 12.7 (for the configuration model).

The overall behavior shown in Fig. 15.4 is typical of percolation in networks. For most degree distributions we expect S to take basically this form with a

⁷A phase transition is *continuous* if the quantity of interest, also called the *order parameter* (S in this case), is zero on one side of the transition and non-zero on the other, but its value is continuous at the transition itself. The alternative to a continuous phase transition is a *discontinuous* or *first-order phase transition*, in which the order parameter jumps discontinuously as it crosses the transition point.

continuous phase transition, as we can demonstrate by the following argument. Suppose the generating function $g_1(u)$ is well-behaved near $u = 1$, having all its derivatives finite,⁸ then we can expand it about this point as

$$\begin{aligned} g_1(u) &= g_1(1) + (u-1)g_1'(1) + \frac{1}{2}(u-1)^2g_1''(1) + O(u-1)^3 \\ &= 1 + \frac{u-1}{\phi_c} + \frac{1}{2}(u-1)^2g_1''(1) + O(u-1)^3, \end{aligned} \quad (15.21)$$

where we have made use of $g_1(1) = 1$ (see Eqs. (12.85) and (15.7)). Substituting into Eq. (15.4), we then find that u satisfies

$$u = 1 + \frac{\phi}{\phi_c}(u-1) + \frac{1}{2}\phi(u-1)^2g_1''(1) + O(u-1)^3 \quad (15.22)$$

or

$$u-1 = \frac{2}{g_1''(1)} \frac{\phi_c - \phi}{\phi_c \phi} + O(u-1)^2. \quad (15.23)$$

We can similarly expand $g_0(u)$ as

$$\begin{aligned} g_0(u) &= g_0(1) + (u-1)g_0'(1) + O(u-1)^2 \\ &= 1 + \frac{2\langle k \rangle}{g_1''(1)} \frac{\phi_c - \phi}{\phi_c \phi} + O(\phi - \phi_c)^2, \end{aligned} \quad (15.24)$$

where we have used $g_0(1) = 1$ and Eqs. (12.87) and (15.23). Substituting into Eq. (15.2) then gives us

$$S = \frac{2\langle k \rangle}{g_1''(1)} \frac{\phi - \phi_c}{\phi_c} + O(\phi - \phi_c)^2. \quad (15.25)$$

In other words, S varies linearly with $\phi - \phi_c$ just above the percolation transition, going to zero continuously as we approach the transition from above. Thus, we would expect the percolation transition to look generically like the curve in Fig. 15.4, with a continuous phase transition as we pass the percolation threshold.⁹

This result is important because it implies that the giant cluster becomes very small as we approach the percolation transition from above. In other words, the network may be “functional” in the sense of having a giant cluster, but the functional portion of the network is vanishingly small. If the network

⁸This excludes the power-law case shown in Fig. 15.3, which is discussed separately later in this section.

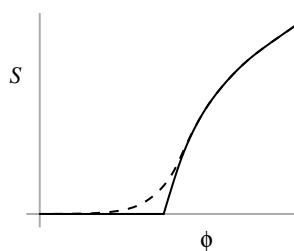
⁹To be more precise, the transition is a *second-order* transition—one where the order parameter is continuous at the transition but its derivative is not.

is a communication network, for example, then a non-zero fraction of all the nodes in the network can communicate with one another so long as there is a giant cluster, but that fraction becomes very small as we approach the percolation threshold, meaning that in practice most nodes are cut off. Thus, one could argue that it is misleading to interpret the percolation threshold as the point where the network stops functioning: in effect most of it has stopped functioning before we reach this point. To fully describe the functional state of the network one should specify not only whether it contains a giant cluster but also what the size of that cluster is.

It is also important to note that the sharp percolation transition of Fig. 15.4 is only truly seen in an infinite network. For networks of finite size—which includes all real networks, of course—the transition gets rounded off. To see this, consider the behavior of the giant cluster in a finite-sized network. Technically, in fact, there *is* no giant cluster for any single finite network. The proper definition of the giant cluster, like the giant component in a random graph, is as a cluster whose size scales in proportion to the size of the network (see Section 11.5). And it makes no sense to talk about the scaling of a cluster with network size when the size of the network is fixed. To get around this problem, let us just consider the largest cluster, which is an acceptable proxy for the giant cluster in a finite-size network. Its size as a fraction of the size of the network should be a reasonable approximation to the size of the giant cluster given by our theory when we are above the percolation transition.

Below the transition the largest cluster will be small in size, but not zero, and hence it fills a small but non-zero fraction of the network, in rough but not perfect agreement with the theoretical prediction $S = 0$. Furthermore, this non-zero value grows as we approach the transition point because clusters in general, including the largest one, grow as the occupation probability ϕ increases. The net result is a slight rounding of the sharp transition predicted by the theory, which is often visible, for example, in computer simulations of percolation on smaller networks. Effects such as this that show up only in finite-sized systems are known as *finite-size effects*.

Even in the limit of large network size there are exceptions to the behavior of Fig. 15.4 and Eq. (15.25). Consider a network with a power-law degree distribution with exponent $2 < \alpha < 3$, as discussed earlier. In this case our assumption that the derivatives of g_1 are finite does not hold (see Fig. 15.3 and the accompanying discussion), so our previous argument breaks down. Not only does the percolation threshold fall at $\phi_c = 0$ for power-law networks, but the giant cluster does not grow linearly as ϕ increases. In general it will grow slower than linearly, the exact functional form depending on the shape of $g_1(u)$



The phase transition at which the giant cluster appears is only sharp in an infinite system (solid line). In a finite-sized system it gets rounded off (dashed line).

near $u = 1$. For example, a typical form is

$$1 - g_1(u) = c(1 - u)^\beta, \quad (15.26)$$

near $u = 1$ with c and β positive constants. Provided $\beta < 1$ this makes the gradient of $g_1(u)$ (and all higher derivatives) infinite at $u = 1$ while still ensuring that $g_1(1) = 1$. With this form for $g_1(u)$, Eq. (15.4) implies that

$$1 - u = (c\phi)^{1/(1-\beta)}. \quad (15.27)$$

Then we can write¹⁰

$$g_0(u) \simeq g_0(1) + g_0'(1)(u - 1) = 1 + \langle k \rangle (u - 1), \quad (15.28)$$

close to $u = 1$, with $\langle k \rangle$ finite so long as the power-law exponent $\alpha > 2$, and hence the giant cluster has size

$$S = \phi[1 - g_0(u)] \simeq \phi \langle k \rangle (1 - u) \sim \phi^{(2-\beta)/(1-\beta)}, \quad (15.29)$$

which goes to zero faster than linearly¹¹ as $\phi \rightarrow 0$ since $(2 - \beta)/(1 - \beta) > 1$ if $\beta < 1$.

Thus we expect the giant cluster to become very small as $\phi \rightarrow 0$. Figure 15.5 shows the equivalent of Fig. 15.4 for a scale-free network with exponent $\alpha = 2.5$, derived from numerical solutions of Eqs. (15.2) and (15.4) and the non-linear form of S close to $\phi = 0$ is clear.

This result mitigates somewhat our earlier statement that scale-free networks are highly robust because $\phi_c = 0$. It is true that the percolation threshold is zero in these networks and hence that there is a giant cluster for any positive ϕ , but that giant cluster can become exceedingly small. A communication network with a power-law degree distribution, for instance, might be formally functional for very small values of ϕ , but in practice the fraction of nodes that could communicate with one another would be so small that the network would probably not be of much use.

¹⁰If we want to be more careful and keep track of the correction terms we can make use of Eq. (12.33) and integrate Eq. (15.26) to show that $g_0(u) = 1 - \langle k \rangle (1 - u) + c(1 - u)^{\beta+1}/(\beta + 1)$. The last term vanishes faster than those before it as $u \rightarrow 1$ because $\beta > 0$ and hence $g_0(u) \simeq 1 - \langle k \rangle (1 - u)$. The form of $g_0(u)$ is at first slightly surprising—one might imagine that the correction term ought to be $O(1 - u)^2$ —but this type of behavior is common with power-law distributions.

¹¹To the extent that one can regard a power-law network as having a percolation transition at $\phi = 0$ it is interesting to ask what the order of this transition is. The answer is unclear since Eq. (15.29) doesn't perfectly fit the standard forms for continuous phase transitions. If we define a transition to be second-order if the order parameter is continuous at the transition and third-order if its derivative is continuous, then the transition is third-order in this case. But one could also argue that the transition is of fractional order between two and three since the order parameter varies as a fractional power of the occupation probability ϕ immediately above the transition.

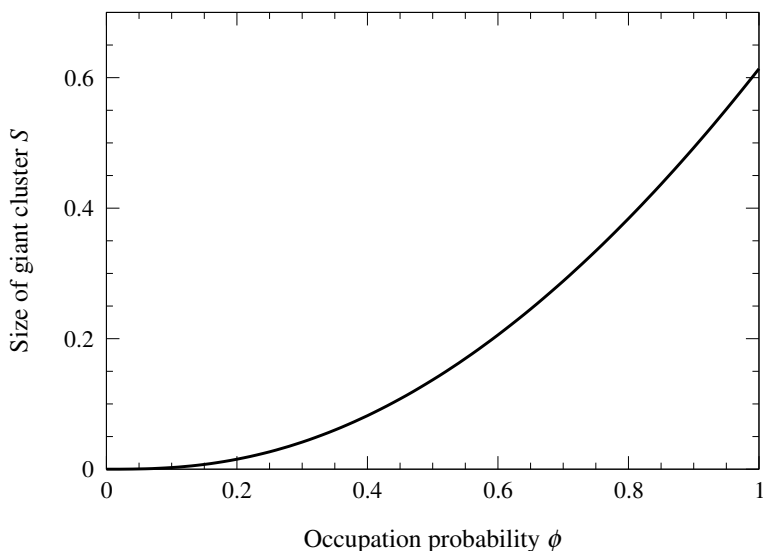


Figure 15.5: Size of the giant cluster for a scale-free network. The size of the giant cluster for a configuration model network with a power-law degree distribution and exponent $\alpha = 2.5$, a typical value for real-world networks. Note the non-linear form of the curve near $\phi = 0$, which means that S , while technically non-zero, becomes very small in this regime. Contrast this figure with Fig. 15.4 for the giant cluster size in a network with an exponential degree distribution.

15.3 NON-UNIFORM REMOVAL OF NODES

In the preceding sections we considered the percolation process in which nodes were removed from a network uniformly at random. This is the classical form of percolation long studied by physicists and mathematicians. When discussing networks, however, it is interesting also to consider other ways in which nodes might be removed. In Section 15.1, for example, we mentioned the possibility of removing nodes in order of their degrees, starting with the highest degrees and working down. This might be effective, for example, as a vaccination strategy for preventing the spread of disease. Should they become infected, the high-degree nodes in a network clearly present a heightened disease risk to the population at large because of their many neighbors, so perhaps vaccinating them first would be a sensible approach.

Let us consider a generalization of our percolation process in which the occupation probability of a node can now depend on its degree. We define ϕ_k to be the probability that a node with degree k is present or occupied

in our network. If ϕ_k is a constant, independent of k , then we recover the uniform scenario of previous sections. Alternatively, if $\phi_k = 1$ for all nodes with degree $k < k_0$ for some constant k_0 , and $\phi_k = 0$ for all $k \geq k_0$, then we effectively remove from the network all nodes with degree k_0 or greater. A host of other choices are also possible, resulting in more complex removal patterns.

Let us again look at percolation on configuration model networks and as before define u to be the average probability that a node is not connected to the giant cluster via one of its neighbors. If the node has degree k , then the probability that it is not connected to the giant cluster via any of its neighbors is u^k and the probability that it is connected to the giant cluster is $1 - u^k$. But in order to belong to the giant cluster, the node itself must also be occupied, which happens with probability ϕ_k , so the probability of being a member of the giant cluster is $\phi_k(1 - u^k)$.

Now we average over the probability distribution p_k of the degree to find the average probability S of being in the giant cluster and get

$$\begin{aligned} S &= \sum_{k=0}^{\infty} p_k \phi_k (1 - u^k) = \sum_{k=0}^{\infty} p_k \phi_k - \sum_{k=0}^{\infty} p_k \phi_k u^k \\ &= f_0(1) - f_0(u), \end{aligned} \quad (15.30)$$

where

$$f_0(z) = \sum_{k=0}^{\infty} p_k \phi_k z^k. \quad (15.31)$$

Note that this new function $f_0(z)$ is not normalized in the conventional manner of a generating function—the value $f_0(1)$ that appears in Eq. (15.30) is not in general equal to 1. Instead it is given by

$$f_0(1) = \sum_{k=0}^{\infty} p_k \phi_k = \bar{\phi}, \quad (15.32)$$

which is the average probability that a node is occupied, or equivalently the average fraction of occupied nodes.

We can calculate the value of u using an approach similar to that for the uniform percolation scenario. The value of u is the probability that you are not connected to the giant cluster via your neighbor, which happens either if the neighbor is not occupied or if it is occupied but it is not connected to the giant cluster via any of its other neighbors. Let k now be the excess degree of the neighboring node. Then the probability that the neighbor is not occupied is $1 - \phi_{k+1}$. Note that the index is $k + 1$ because ϕ_k is defined in terms of the total degree of a node, which is one greater than the excess degree (see

Section 12.2). The probability that the neighbor *is* occupied but is itself not connected to the giant cluster is $\phi_{k+1}u^k$. Adding up the terms and averaging over the distribution q_k of the excess degree, we then find that

$$u = \sum_{k=0}^{\infty} q_k(1 - \phi_{k+1} + \phi_{k+1}u^k) = 1 - f_1(1) + f_1(u), \quad (15.33)$$

where

$$f_1(z) = \sum_{k=0}^{\infty} q_k \phi_{k+1} z^k \quad (15.34)$$

and we have used $\sum_k q_k = 1$.

The definition of $f_1(z)$ looks slightly odd because of the subscript $k + 1$ on ϕ_{k+1} . If we prefer, we can write it using the full expression for the excess degree distribution, Eq. (15.3), which gives

$$f_1(z) = \frac{1}{\langle k \rangle} \sum_{k=0}^{\infty} (k+1)p_{k+1}\phi_{k+1}z^k = \frac{1}{\langle k \rangle} \sum_{k=1}^{\infty} kp_k\phi_k z^{k-1}, \quad (15.35)$$

which has a more symmetric look to it. We can also write

$$f_1(z) = \frac{f_0'(z)}{g_0'(1)}, \quad (15.36)$$

where $g_0(z)$ is defined as before. This expression can be useful for calculating $f_1(z)$ once $f_0(z)$ has been found. Note also that, like $f_0(z)$, the function $f_1(z)$ is not normalized in the conventional manner of a generating function—the value of $f_1(1)$ is not in general equal to 1.

Equations (15.30) and (15.33), which were first given by Callaway *et al.* [93], give us a complete solution for the size of the giant cluster for our generalized percolation process. As an example of their use, consider again a network with exponential degree distribution given by Eq. (15.12) and suppose we remove all nodes that have degree k_0 or greater. That is, we choose

$$\phi_k = \begin{cases} 0 & \text{if } k \geq k_0, \\ 1 & \text{otherwise.} \end{cases} \quad (15.37)$$

Then we have

$$f_0(z) = (1 - a) \sum_{k=0}^{k_0-1} (az)^k = [1 - (az)^{k_0}] \frac{1 - a}{1 - az}, \quad (15.38)$$

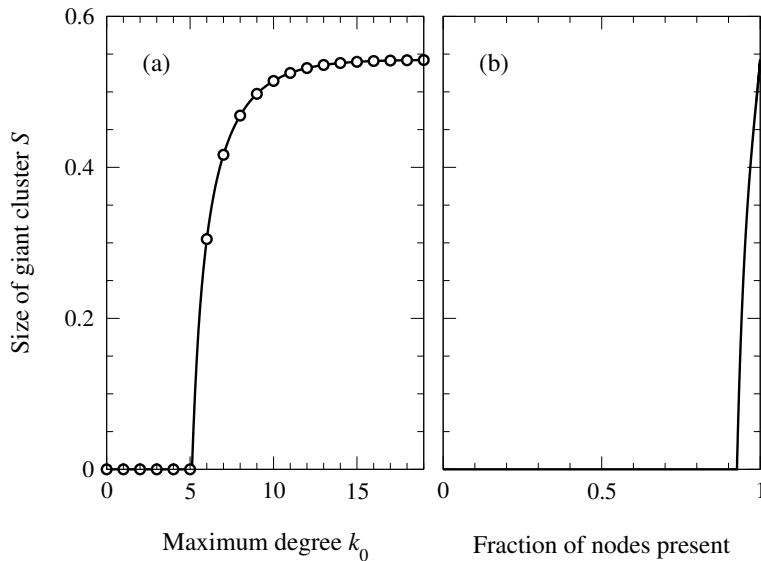


Figure 15.6: Size of the giant percolation cluster as the highest degree nodes in a network are removed. (a) The size of the giant cluster in a network with an exponential degree distribution $p_k \sim a^k$ with $a = 0.6$ as nodes are removed in order of degree, starting from those with the highest degree. The curve is shown as a function of the degree k_0 of the highest-degree node remaining in the network. Technically, since k_0 must be an integer, the plot is only valid at the integer points marked by the circles; the curve is just an aid to the eye. (b) The same data plotted now as a function of the fraction $\bar{\phi}$ of nodes remaining in the network.

and

$$f_1(z) = \frac{f'_0(z)}{g'_0(1)} = \left[(1 - (az)^{k_0}) - k_0(az)^{k_0-1}(1 - az) \right] \left(\frac{1 - a}{1 - az} \right)^2. \quad (15.39)$$

For this choice Eq. (15.33) becomes a polynomial equation of order k_0 and unfortunately such equations are not solvable exactly for their roots (unless $k_0 \leq 4$). It is, however, fairly easy to find the roots numerically, particularly given that we know that the root of interest in this case lies in the range between zero and one, then we can calculate the size of the giant cluster from (15.30).

Figure 15.6a shows the results of such a calculation, plotted as a function of k_0 . Looking at this figure, consider what happens as we lower k_0 from an initial high value, effectively removing more and more of the high-degree nodes in our network. As the figure shows, the size of the giant cluster decreases only slowly at first. This is because there are not many nodes of very high degree in

the network, so very few have been removed. Once k_0 passes a value around 10, however, our attack on the network starts to become evident in a shrinking of the giant cluster, which becomes progressively more rapid until the size of the cluster reaches zero around $k_0 = 5$.

One might be forgiven for thinking that Fig. 15.6a portrays a network quite resilient to the removal of even its highest-degree nodes. We have to remove nodes all the way down to degree five in order to break up the giant cluster. This impression is misleading, however, because it fails to take account of the fact that the vast majority of nodes in the network are of low degree, so that even when we have removed all nodes with degree greater than 5, we have still removed only a small fraction of all nodes.

Perhaps a more useful representation of the situation is to make a plot of the giant cluster size as a function of the fraction $\bar{\phi}$ of occupied nodes in the network, which is

$$\bar{\phi} = f_0(1) = 1 - a^{k_0}. \quad (15.40)$$

Figure 15.6b shows the result replotted in this way and reveals that the giant cluster in fact disappears completely when only about 8% of the highest-degree nodes in the network have been removed. By contrast, when we removed nodes uniformly at random, as shown in Fig. 15.4, we had to remove nearly 70% of the nodes to destroy the giant cluster. Though the difference is startling, however, it is also intuitively reasonable. The high-degree nodes have a lot of connections, all of which are lost if we remove those nodes.

These results suggest, for example, that were we able to find the highest degree nodes in a network of disease-causing contacts and vaccinate them, thereby effectively removing them from the network, it would be a much more efficient strategy for disease control than simply vaccinating at random.

A particularly striking example of the effect described here arises in networks with power-law degree distributions. In these networks, as we have seen, uniform removal of nodes *never* destroys the giant cluster, provided the exponent of the power-law lies between 2 and 3. In contrast, removal of the highest-degree nodes in these networks has a devastating effect. Once again we cannot solve for S in closed form but it is straightforward to perform a numerical solution. Figure 15.7a shows the equivalent of Fig. 15.6b for the power-law case and as we can see the giant cluster disappears extraordinarily rapidly as the high-degree nodes are removed. Only a few percent of the nodes need be removed to completely destroy the giant cluster, the exact fraction depending on the exponent of the power law.

If we want to calculate this fraction we can do so by observing that the phase transition at which the giant cluster appears or disappears falls at the

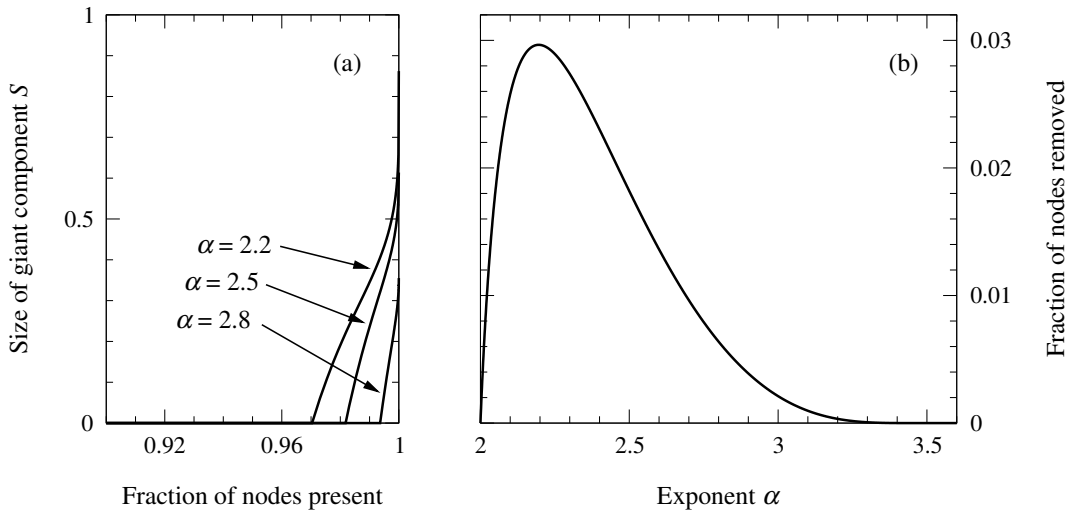


Figure 15.7: Removal of the highest-degree nodes in a scale-free network. (a) The size of the giant cluster in a configuration model network with a power-law degree distribution as nodes are removed in order of their degree, starting with the highest-degree nodes. Only a small fraction of the nodes need be removed to destroy the giant cluster completely. (b) The fraction of nodes that must be removed to destroy the giant cluster as a function of the exponent α of the power-law distribution. For no value of α does the fraction required exceed 3%.

point where the non-trivial solution of Eq. (15.33) appears or disappears, which is the point at which the right-hand side of the equation is tangent to the line $y = u$ at $u = 1$. That is, the general criterion for the transition point is

$$f_1'(1) = 1. \quad (15.41)$$

(Alternatively, we could say that the giant cluster exists if and only if $f_1'(1) > 1$.) Again, exact solutions of (15.41) are often not possible but we can solve numerically. Doing this for the power-law case we find the results shown in Figure 15.7b, which plots the fraction of nodes that need be removed to destroy the giant cluster as a function of the exponent α . As we can see, the curve peaks around $\alpha = 2.2$ at a value just below 3%. Thus in no case need we remove more than 3% of nodes to destroy the connectivity in the network.

Scale-free networks are thus paradoxically both robust and fragile, a point first emphasized by Albert *et al.* [17]. On the one hand, they are remarkably robust to the random failure of their nodes, with the giant cluster persisting no matter how many nodes we remove. (Although one should bear in mind the *proviso* of Section 15.2.1 that the size of the giant cluster matters also, and this becomes very small when the fraction ϕ of occupied nodes tends to zero.)

On the other hand, scale-free networks are very fragile to attacks targeted specifically at their highest-degree nodes. We need remove only the tiniest fraction of the high-degree hubs in such a network to entirely destroy the giant cluster.

The fragility of scale-free networks to targeted attacks is both bad news and good news. It is bad news for networks such as the Internet that we wish to defend against possible attack: a communication network that can easily be brought down by a malicious adversary targeting just a few of its most crucial nodes may be a disaster waiting to happen.

On the other hand, results like these could also be exploited to help eradicate or reduce disease by targeting vaccination efforts at network hubs. It is worth noting, however, that it's not necessarily easy to find the hubs in a network, so that implementation of a targeted vaccination strategy may be difficult. In most cases one does not know the entire network and so cannot simply pick out the high-degree nodes from a list.

One intriguing way of getting around this problem has been put forward by Cohen *et al.* [115], who suggest that we make use of the structure of the network itself to find the high-degree nodes. In their scheme, which they call "acquaintance immunization," they propose that one choose members of the population at random and then get each of them to nominate an acquaintance. Then that acquaintance receives a vaccination against the disease under consideration. The acquaintance in this scenario is a "node at the end of an edge," so in the configuration model this node would have degree distributed according to the excess degree distribution, Eq. (12.16), rather than the original degree distribution of the network. But the excess degree distribution, as discussed in Section 12.2, is biased towards high-degree nodes since there are more edges that end at a high-degree node than at a low-degree one. Thus the selection of individuals in the scheme of Cohen *et al.* is also biased towards those with high degree. The selected individuals are not guaranteed to be the highest-degree nodes in the network, but we are a lot more likely to find the hubs this way than if we just choose nodes at random, and in simulations the acquaintance immunization scheme appears to work quite well.

On the other hand it has some drawbacks too. First, contact networks in the real world are of course not configuration models and it is unclear how accurately the theoretical results describe real situations. Second, real contact networks mostly don't have power-law degree distributions, instead having somewhat shorter tails than the typical power law, which will reduce the effectiveness of the scheme, or indeed of any scheme based on targeting the highly connected nodes. Another issue is that, in asking people to name their acquaintances, the acquaintance immunization scheme necessarily probes the

network of who is acquainted with whom, which is in general not the same as the network of disease transmission. We can do our best to make the networks similar, asking participants to name only acquaintances they have seen recently and in person, rather than those they might not have seen for a while or might only have talked to on the phone. Still, the difference between the two networks means that the scheme might end up focusing vaccination efforts on some of the wrong people. Finally, one cannot guarantee that the acquaintances identified by the process can be located, or that they will consent to being vaccinated. Possibly a ticket-based scheme for encouraging participation akin to the “respondent-driven sampling” method of Section 4.7 could be used to improve the success rate.

15.4 PERCOLATION IN REAL-WORLD NETWORKS

Having seen how percolation plays out in model networks, let us now take a look at some real ones. If we have data on the structure of a network, then we can simulate the percolation process on a computer, removing nodes one by one and examining the resulting clusters. Although this is straightforward in theory, it requires some care to get good results in practice. The main issue is that the percolation process is a random one: the nodes are removed in random order, which means that the cluster sizes can vary depending on the precise order we choose. Even in the case where nodes are removed in order of their degree the process is still random to some extent, since there can be many nodes with a given degree, among which, to avoid bias, we typically choose at random.

Randomness can easily be simulated on a computer using standard random number generators, but the results of the simulation will then vary from one run to another, depending on the specific random numbers generated. So to get a reliable picture of how percolation affects a network we must perform the entire calculation many times, removing the nodes in a different random order each time, so that we can see what the typical behavior is, as well as the range of variation from run to run. This in turn means that we need to be able to perform the percolation calculation quickly. In a typical situation we might want to repeat the percolation calculation a thousand times with different random orders of removal and even if each calculation took just one minute of computer time, a thousand runs would still take a day.

If we are crafty, however, we can do much better than this and get an answer in just a few seconds for networks of the typical sizes we have been considering in this book.

15.5 COMPUTER ALGORITHMS FOR PERCOLATION

The actual process of percolation itself—the random occupation or removal of nodes in a network—is trivial to simulate on a computer. The more challenging task is finding the percolation clusters this process creates, which are the primary object of interest in percolation calculations. The most straightforward way to find the clusters is to make use of the breadth-first search algorithm of Section 8.5.4. Recall that this algorithm can find all components in a network in time $O(m + n)$, where m is the total number of edges in the network and n is the total number of nodes, or just $O(n)$ for a sparse network in which $m \propto n$. If we remove a set of nodes from a network, along with the edges attached to them, then the resulting percolation clusters are by definition the components of the network that remains, and hence we can use breadth-first search to find clusters as well.

In the case of uniform random removal of nodes, for instance, a numerical calculation of the size of the largest percolation cluster for a given network would involve going through each node in turn and removing it (and its edges) from the network with probability $1 - \phi$, then finding the resulting clusters and examining them to find the size S of the largest one. Then we would repeat the entire process, starting with the complete network again, removing a different randomly chosen set of nodes, and finding the clusters. Repeating the calculation a large number of times, we can calculate a mean value $S(\phi)$ for the size of the largest cluster when nodes are present or functioning with probability ϕ .

If we are interested in only a single value of ϕ , this is, in fact, the best algorithm to use and the fastest known way of getting an answer. Usually, however, we are interested, as in previous sections, in the behavior of the system over the whole range of ϕ from zero to one, or at least some portion of that range. In that case, we would have to repeat the entire calculation above for many values of ϕ in the range of interest. This process is time-consuming and is not the best way to approach the problem.

Consider instead the following alternative approach, which appears at first to be only a slight variation on the previous one, but leads, as we will see, to much more efficient algorithms. Instead of making each node in the network occupied with independent probability ϕ , let us make a fixed number r of nodes occupied, repeating the calculation many times for a given value of r and averaging to get a figure S_r for the size of the largest cluster (or any other quantity of interest) as a function of r .

The calculation doesn't directly give us the result we want: S_r is not the same as $S(\phi)$ and it is the latter we are interested in. If, however, we know the

value of S_r for every value of the integer r from 0 to n , then we can calculate $S(\phi)$ as follows. If each node in the network is occupied with probability ϕ , then the probability that there are exactly r nodes occupied is given by the binomial distribution

$$p_r = \binom{n}{r} \phi^r (1 - \phi)^{n-r}. \quad (15.42)$$

Averaging over this distribution, the average size of the largest cluster as a function of ϕ is then

$$S(\phi) = \sum_{r=0}^n p_r S_r = \sum_{r=0}^n \binom{n}{r} \phi^r (1 - \phi)^{n-r} S_r. \quad (15.43)$$

At first sight, this appears to be a less promising approach for calculating $S(\phi)$ than the previous one. To make use of Eq. (15.43) we need to know S_r for all r , and the calculation takes time $O(m + n)$ for one value of r using breadth-first search, so it is going to take $O(n(m + n))$ to calculate all n values, or $O(n^2)$ if $m \propto n$. Given that we also need to repeat the calculations many times to average over the randomness, the entire process could take a very long time to complete.

There is, however, a faster way to calculate S_r for all r , inspired by the simple observation that if we have already found all the clusters in a network with r nodes present, then we can find the clusters with $r + 1$ nodes simply by adding one more node. Most of the clusters will not change when we add just one node and by concentrating only on those that do change we can save ourselves the work of performing an entire new breadth-first search, and hence save ourselves a lot of computer time. A simple algorithm for doing this is as follows.

Rather than removing nodes from the complete network, this algorithm works by building the network up from an initial state in which no nodes are occupied and occupying nodes one by one until we recover the entire network. As we add each node to the network we also add the accompanying edges that join it to other nodes. Only connections to other nodes that are already present need be added.

For the purposes of the algorithm, let us break down this process as shown in Fig. 15.8. Each new node is first added with, initially, no accompanying edges (Fig. 15.8a). In this state it forms a cluster all on its own. Then, one by one, we add its edges, those that connect it to other nodes already present. If there are no such edges then our new node remains a cluster on its own. If there are edges to be added, however, then the first edge we add necessarily joins our node to an adjacent cluster—see Fig. 15.8b. Subsequent edges are more complicated. They can do one of two things. An edge can connect our node

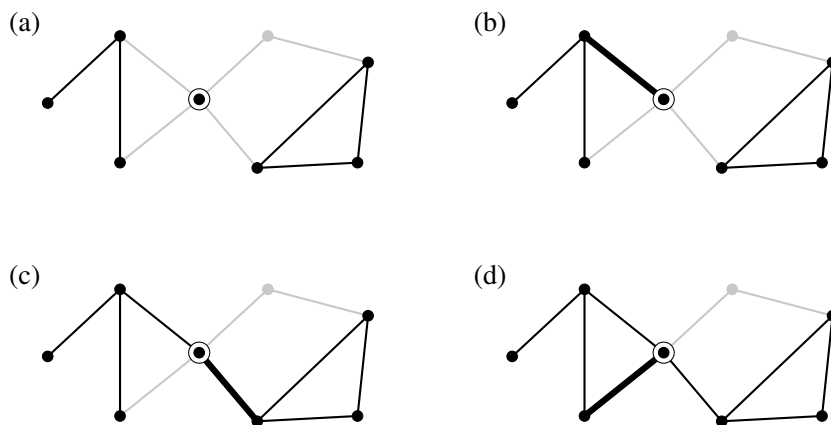


Figure 15.8: Computer algorithm for percolation. In the percolation algorithm described in the text we add nodes to our network one by one, rather than taking them away. Each addition consists of several steps. (a) We add the node itself but none of its accompanying edges yet. At this stage the node constitutes a new cluster in its own right. (b) We start adding the accompanying edges (if any) in any order we like. Only edges that connect to other nodes already present in the network are added. The first edge added (if any are added) will thus, by definition, always join the new node to one of the previously existing clusters. Or to put it another way, it will join two clusters together, one of the old clusters and the new cluster that consists of just the single added node. (c) In this example the next edge added also joins two clusters together. (d) The final edge added joins two nodes that are already members of the same cluster, so the cluster structure of the network does not change.

to another, different cluster, in which case in the process it joins two clusters together, making them into a single cluster—see Fig. 15.8c. Alternatively, it could join our node to another member of the same cluster it already belongs to, as in Fig. 15.8d. In this case, no clusters are joined together, and in terms of the size and identity of the clusters in the network the added edge has no effect.

To keep track of the clusters in the network, therefore, our algorithm needs to do two things. First, when an edge is added it needs to identify the clusters to which the nodes at either end belong. Second, if the clusters are different, it needs to join them together into a single cluster. (If they are the same nothing need be done.)

There are various ways of achieving this but a simple one is just to put a label, such as an integer, on each node, denoting the cluster to which it belongs—see Fig. 15.9a. Then it is a simple matter to determine whether two nodes belong to the same cluster—they do if their labels are the same—and joining two clusters together is just a matter of relabeling all the nodes in one of the clusters to

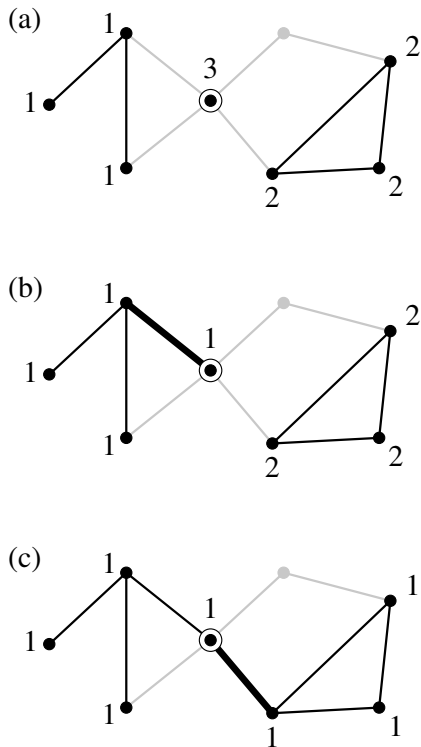


Figure 15.9: Using labels to keep track of clusters. In the algorithm described in the text, each node is given a label, typically an integer, to denote which cluster it belongs to. In this example there are initially two clusters, labeled 1 and 2, and a new node is added between them. (a) The new node is initially added without its accompanying edges and forms a new cluster, which is labeled number 3. (b) An edge is added that connects cluster 3 to cluster 1, so we relabel one cluster to give it the same label as the other. In the algorithm described in the text we always relabel the smaller of the two clusters, which is cluster 3 in this case. (c) The next edge added joins clusters 1 and 2 and we relabel cluster 2.

match the label of the other cluster. This process is illustrated in Fig. 15.9.

The full algorithm is as follows:

1. Start with an empty network with no occupied nodes. Let $c = 0$ be the number of clusters in the network initially. Choose at random an order in which the nodes will be added to the network.
2. Add the next node in the chosen order, initially with no edges. This node forms a cluster in its own right, so increase c by one and label the node with label c to indicate which cluster it belongs to. Also make a note that cluster c has size 1.
3. Go through the edges attached to this node one by one in any order. For each edge determine whether the node at the other end has already been added to the network. If it has, add the edge to the network.
4. For each edge added, examine the cluster labels of the nodes at either end. If they are the same, do nothing. If they are different, choose one of the clusters and relabel all its nodes to have the same label as the other cluster. Update the record of the size of the cluster to be equal to the sum

of the sizes of the two clusters from which it was formed.

5. Repeat from step 2 until all nodes have been added.

At the end of this process, we have gone from an entirely empty network to the complete final network with all nodes and edges present and in between we have passed through a state with every possible intermediate number r of nodes. Moreover, in each of those states we had a complete record of the identities and sizes of all the clusters, which we can use, for instance, to find the size S_r of the largest cluster. Then we can feed the results into Eq. (15.43) to calculate $S(\phi)$ for any ϕ . As discussed previously, we typically want to average the results over many runs of the algorithm to allow for random variations from one run to another, which arise from variations in the order in which the nodes are added. This, however, is no longer a serious impediment to finishing the calculation because, if implemented appropriately, the algorithm can be made to run very quickly.

The most time-consuming part of the algorithm is the relabeling of clusters when they are joined together. Note, however, that when an edge joins two different clusters we are free to choose which of the two we relabel. It turns out that the speed of the algorithm can be improved greatly if we choose always to relabel the smaller one. (If the two clusters have the same size, it does not matter which we choose to relabel.) To see this, consider the following argument.

If we always relabel the smaller of two clusters, then the relabeled cluster must have been joined with one at least as large as itself and hence it is now a part of a cluster at least twice its previous size. Thus, every time a node is relabeled, the cluster it belongs to at least doubles in size. Given that each node starts off as a cluster of size 1, the size of the cluster to which it belongs after k relabelings is at least 2^k . And since no node can belong to a cluster of size greater than the size n of the whole network, the maximum number of relabelings a node can experience during the entire algorithm is given by $2^k = n$ or $k = \log_2 n$, and the maximum number of relabeling operations on all n nodes is $n \log_2 n$. In practice, the relabeling itself can be done using breadth-first search, following edges from each relabeled node to its neighbors to see whether they require relabeling too. The average number of edges followed from a node is equal to the average node degree $2m/n$ (Eq. (6.15)), so the time taken per node is $O(1 + m/n)$, including the (constant) time required to update the label itself. Thus, the total time to perform the relabeling part of the algorithm is at most $O(1 + m/n) \times n \log n = O((m + n) \log n)$.

The other parts of the algorithm are typically faster than this. The adding of the nodes takes $O(n)$ time and the adding of the edges takes $O(m)$ time. So the overall running time of the algorithm to leading order is $O((m + n) \log n)$,

or $O(n \log n)$ on a sparse network with $m \propto n$, which is much better than the $O(n(m + n))$ for our first naive algorithm.

The same algorithm can also be used when nodes are added or removed non-uniformly. For instance, if nodes are to be removed in decreasing order of their degrees we simply reverse that process and add nodes to an initially empty network in increasing order of degrees. The details of the algorithm itself are unchanged—only the order of the nodes changes.

This algorithm works well in practice for almost all calculations. It is not, however, the very fastest algorithm for the percolation problem. There exists an even faster one, which runs in $O(m + n)$ time (or $O(n)$ for a sparse network) and is also considerably simpler to program, although its outward simplicity hides some subtleties. The reader interested in learning more about this approach is encouraged to look at Ref. [371].

15.5.1 RESULTS FOR REAL-WORLD NETWORKS

Figure 15.10 shows results of percolation calculations on four different networks using the algorithm of the previous section. In these calculations the occupied nodes were chosen uniformly at random and the figure shows in each case the size S of the largest cluster as a fraction of system size, plotted as a function of the fraction ϕ of occupied nodes. As described in Section 15.2.1, the largest cluster acts as a proxy for the giant cluster in numerical calculations on fixed networks for which the idea of a giant cluster, as a cluster that scales with system size, is meaningless.

The top two networks in the figure, a power grid and a road network, are both networks with non-power-law degree distributions—the power grid has a roughly exponential distribution while the road network has only nodes of degree one to four and nothing else. For these networks we expect to see behavior of the generic type described in Section 15.2.1: a continuous percolation transition at a non-zero value of ϕ from a regime in which $S \simeq 0$ to a regime of non-zero S . Because the networks are relatively small, however (4941 nodes for the power grid, 935 for the road network), we also expect to see some rounding of the transition (see Section 15.2.1).

And this is in fact what we do see. In these two cases S is close to zero below a certain value of ϕ , then grows rapidly but with a certain amount of rounding near the transition. Overall, other than the rounding, the shape of the curves is qualitatively similar to that of Fig. 15.4. One could even tentatively make an estimate of the position of the percolation threshold, which appears to fall around $\phi = 0.6$ or 0.7 in both networks.

The bottom two frames in the figure tell a different story. These show results

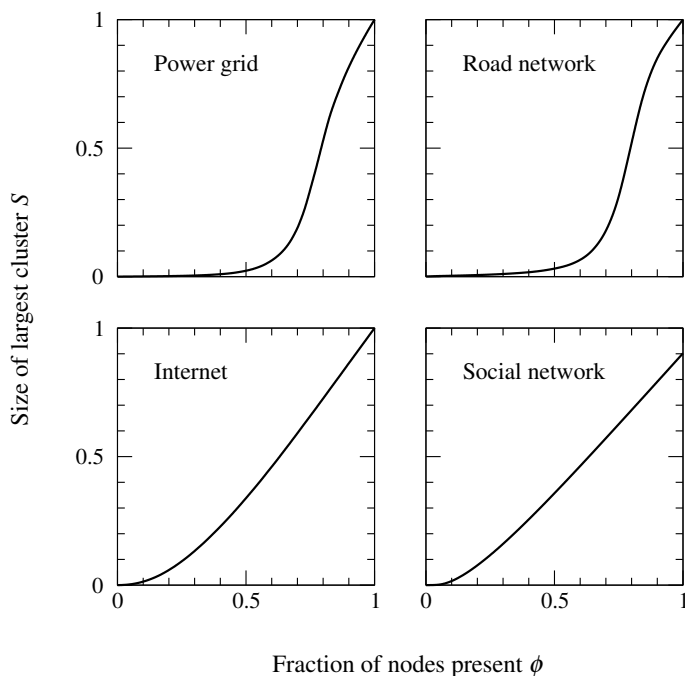


Figure 15.10: Size of the largest cluster as a function of occupation probability for percolation on four networks. The four frames of this figure show the size of the largest percolation cluster, measured as a fraction of network size, for random removal of nodes from four real-world networks: the western United States power grid, the network formed by the US interstate highways, the Internet at the level of autonomous systems, and a social network of professional collaborations between physicists. Each curve is averaged over 1000 random repetitions of the calculation, which is why the curves appear smooth.

for percolation on the Internet and a social network. Both of these networks have approximately power-law degree distributions and thus, based on the insights of Section 15.2.1, might be expected to show no percolation transition (or a transition at $\phi = 0$ if you prefer) and non-linear growth of the largest cluster with growing ϕ . Again our expectations seem to be borne out, at least qualitatively, by the numerical results. In both networks the value of S appears to take non-zero values for all $\phi > 0$ and the initial growth for small ϕ shows some curvature, indicating non-linear behavior.

Thus our percolation theory for the configuration model seems in this case to provide a good general guide to the behavior of real-world networks. The

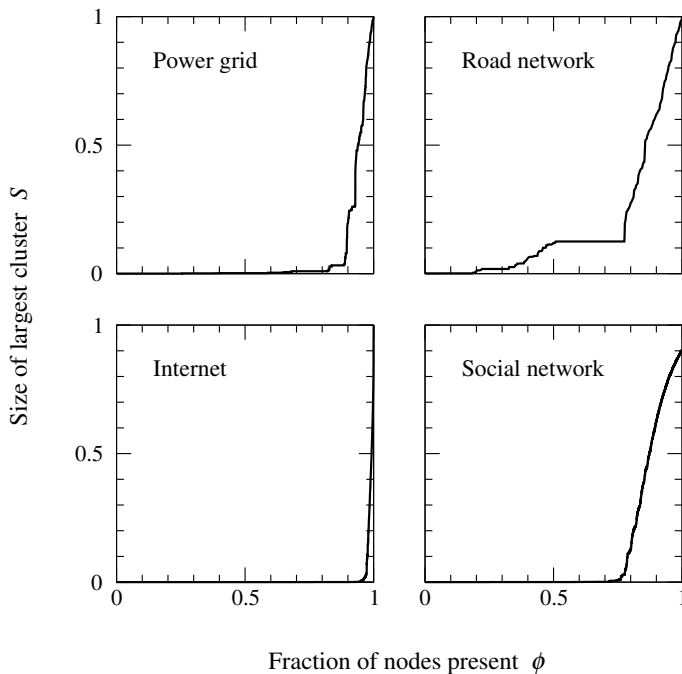


Figure 15.11: Size of the largest percolation cluster as a function of occupation probability for targeted attacks on four networks. The four frames in this figure show the size of the largest cluster, as a fraction of network size, when nodes are removed in degree order, highest degrees first, from the same four networks as Fig. 15.10. Since this is mostly a deterministic process and not a random one (except for random choices between nodes of the same degree) the curves cannot be averaged as in Fig. 15.10 and so are relatively jagged.

power-law networks are robust against random removal of nodes, in the sense that a giant cluster remains even when most nodes have been removed. The non-power-law networks, by contrast, become essentially disconnected after relatively few nodes have been removed—just about 40% in this case.

Figure 15.11 shows results for the same four networks when nodes are removed in order of degree, highest degrees first. As we can see, this “attack” on the network is more effective at reducing the size of the largest cluster than is random removal, for all four networks. However, the difference between Figs. 15.10 and 15.11 is not so great for the first two networks, the power grid and the road network. The giant cluster in both of these networks survives nearly as long under the targeted attack as under random removal. This is as

we would expect, since neither has a significant number of very high-degree nodes (the road network, with maximum degree four, has none at all), so that removal of the highest-degree nodes is not so very different from removal of randomly chosen nodes.

For the second two networks, however, the Internet and the social network, which both have roughly power-law degree distributions, the effect is far larger. Where these networks were more resilient to random removal than the others, they are clearly less resilient, at least by this measure, to targeted attack. The Internet in particular has a largest cluster size that falls essentially to zero when only about 5% of its highest-degree nodes have been removed, a behavior similar again to our theoretical calculations (see Fig. 15.7 on page 591). Thus, the real Internet appears to show the mix of robust and fragile behavior that we saw in our calculations for the configuration model, being remarkably resilient to the random removal of nodes but far more susceptible to targeted attacks.

Overall, therefore, percolation theory seems to be successful as a qualitative guide to the robustness of networks. Its predictions are not perfectly accurate in general, but it gives a good feel for what we should expect to see as nodes fail or are removed.

In the next chapter we will see another application of percolation, to the spread of diseases on networks.

EXERCISES

15.1 Consider a site percolation process in which nodes are removed uniformly at random from a random 4-regular network (i.e., a configuration model where all nodes have degree 4). You can assume the network is large.

- a) Give an expression for the size S of the giant percolation cluster as a fraction of total network size.
- b) Find the critical occupation probability ϕ_c .
- c) Find the value of ϕ at which $S = 1$. This implies that the giant cluster fills the whole network. How can this happen, given that the most it can fill is the whole of the giant component?

15.2 You may find it useful to solve Exercise 12.13 first before solving this one.

Consider the site percolation process with uniform random occupation of nodes and occupation probability ϕ on a Poisson random graph with mean degree c .

- a) Show that the network formed by the occupied nodes and their edges is itself a Poisson random graph. What is the mean degree of a node in this graph?

- b) Hence show that the mean size of a small cluster in the non-percolating regime (no giant cluster) is

$$\langle s \rangle = \frac{\phi}{1 - c\phi}.$$

- c) Using the results of Section 12.10.9 show that the probability π_s that a randomly chosen node belongs to a cluster of size s is

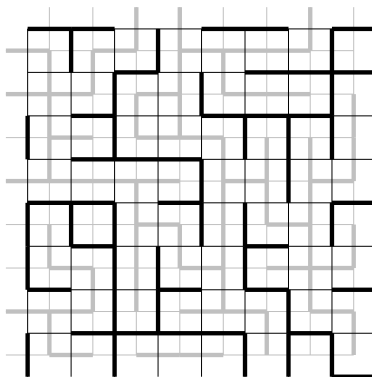
$$\pi_s = \begin{cases} 1 - \phi & \text{if } s = 0, \\ \phi e^{-sc\phi} (sc\phi)^{s-1}/s! & \text{if } s > 0. \end{cases}$$

15.3 Consider a configuration model network with nodes of degree 1, 2, and 3 only, in fractions p_1 , p_2 , and p_3 , respectively.

- Find the value of the critical node occupation probability ϕ_c at which site percolation takes place on the network.
- Show that there is no giant cluster for any value of the occupation probability ϕ if $p_1 > 3p_3$. Why does this result not depend on p_2 ?
- Find the size of the giant cluster as a function of ϕ . (Hint: you may find it useful to remember that $u = 1$ is always a solution of the equation $u = 1 - \phi + \phi g_1(u)$.)

15.4 Use Eq. (15.7) to calculate the position of the percolation threshold for uniform random removal of nodes from a configuration model network with the (properly normalized) exponential degree distribution $p_k = (1 - a)a^k$ with $a < 1$. Verify that the answer you get is the same as that given (by a different method) in Eq. (15.20).

15.5 Consider the problem of (uniform) bond percolation on a square lattice and consider the following construction:



Here we have taken a bond percolation system (in black) and constructed another one interlocking it (in gray), such that a bond in the new system is occupied if and only if the intersecting bond on the old system was not. Such an interlocking system is called a *dual lattice*.

- If the fraction of occupied bonds on the original lattice is ϕ , what is the fraction of occupied bonds on the dual lattice?

- b) Show that there is a path from top to bottom of the dual lattice if and only if there is no path from side to side of the original lattice.
- c) Hence argue that the percolation transition on an infinite square lattice occurs at $\phi = \frac{1}{2}$.

15.6 Consider a (uniform) bond percolation process with edge occupation probability ϕ on a random graph with Poisson degree distribution and mean degree c , in the limit of large network size n .

- a) Write down an equation whose solution gives the probability u that a node is not connected to the giant percolation cluster via a particular one of its edges.
- b) In terms of u write down an expression for the probability that a node is not in the giant cluster given that it has degree k .
- c) Hence, or otherwise, write down an expression in terms of u , c , and k for the probability that a node has degree k given that it is not in the giant cluster.
- d) Thus, show that the mean degree of nodes not in the giant cluster is cu .

15.7 In Section 15.3 we examined what happens when the highest-degree nodes are removed from a configuration model network with a power-law degree distribution.

- a) For the case of the “pure” power-law degree distribution $p_k = k^{-\alpha}/\zeta(\alpha)$ for $k \geq 1$ (and $p_0 = 0$), show that the phase transition at which the giant cluster disappears occurs when all nodes with degree $k > k_0$ have been removed, where k_0 satisfies

$$\sum_{k=1}^{k_0} (k^{-\alpha+2} - k^{-\alpha+1}) = \zeta(\alpha - 1).$$

- b) Using the fact that $\sum_{k=1}^{k_0} k^{-x} + \sum_{k=k_0+1}^{\infty} k^{-x} = \zeta(x)$, and making use of the trapezoidal rule (Eq. (13.103) on page 469) for large values of k , show that

$$\sum_{k=1}^{k_0} k^{-x} \simeq \zeta(x) - \frac{1}{2}(k_0 + 1)^{-x} - \frac{(k_0 + 1)^{-x+1}}{x - 1}.$$

- c) Keeping leading-order terms in k_0 only, show that the giant cluster disappears approximately at the point where

$$(k_0 + 1)^{-\alpha+3} = (\alpha - 3)[\zeta(\alpha - 2) - 2\zeta(\alpha - 1)].$$

- d) Find the approximate value of k_0 at the point where the giant cluster disappears for $\alpha = 2.5$.

15.8 Consider the problem of non-uniform percolation on a configuration model network, as discussed in Section 15.3, where the occupation probability ϕ_k for a node is a function of degree k .

- a) Show, by a graphical argument or otherwise, that a giant cluster can exist in the network only if $f_1'(1) > 1$, where $f_1(z)$ is the function defined in Eq. (15.34).

- b) A configuration model network has a (properly normalized) exponential degree distribution $p_k = (1-a)a^k$ with $a < 1$ and an occupation probability $\phi_k = b^k$ with $b < 1$, so that high-degree nodes are more likely to be removed than low-degree ones. Show that the system has a giant cluster if $2ab^2(1-a)^2 > (1-ab)^3$.

15.9 Recall the “acquaintance immunization” scheme discussed in Section 15.3. In this scheme, instead of vaccinating people at random, we ask people to nominate a friend then we vaccinate the friend. Because one’s friends tend to be the popular people, this has the beneficial effect of vaccinating people with many contacts.

Consider an acquaintance immunization process on an arbitrary network, in which a fraction f of the population are chosen uniformly at random and then each of them nominates one of their friends uniformly at random to receive the vaccination.

- a) Show that the expected number of nominations received by individual i is $f\kappa_i$, where κ_i is the sum of the reciprocals of the degrees of i ’s neighbors in the network:

$$\kappa_i = \sum_j \frac{A_{ij}}{k_j}.$$

- b) Hence argue that the probability that an individual is not vaccinated, i.e., the probability that the corresponding network node is occupied in the sense of percolation, is $e^{-f\kappa_i}$ when n is large.
- c) Between two nodes with the same degree, of which one has high-degree neighbors and one has low-degree neighbors, which is more likely to be vaccinated? Is this a good thing or a bad thing as far as preventing disease is concerned?

15.10 Consider the computer algorithm for percolation described in Section 15.5, but suppose that upon the addition of an edge between two clusters we relabel not the smaller of the two clusters but one or the other cluster chosen at random. Show by an argument analogous to the one in Section 15.5 that the worst-case running time of this algorithm is $O(n^2)$ —substantially worse than the $O(n \log n)$ of the algorithm that always relabels the smaller cluster.

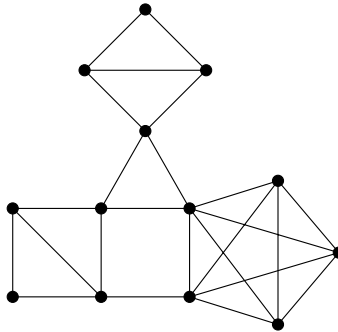
15.11 Suppose we know the contact network over which a disease is spreading and we are given a certain “budget” for vaccinating nodes, meaning we can vaccinate a specified number r of nodes. Our goal is to minimize the number of people who catch the disease.

- a) Suppose that after removing our r nodes the network consists of a set of k percolation clusters with sizes $s_1 \dots s_k$. If the disease starts at a single randomly chosen node in the network and spreads to all nodes in the same percolation cluster as that node (or doesn’t spread if the node itself has been vaccinated), show that the expected number I of nodes that get infected is

$$I = \frac{1}{n} \sum_{m=1}^k s_m^2.$$

Hence argue that the smallest possible value for the expected size of a disease outbreak is $I = (n - r)^2 / nk$ when there are k clusters. All other things being equal, therefore, it's better to allocate the vaccinations so as to divide the network into the largest number of clusters possible. In effect, one should use the vaccinated nodes to create "firewalls" that divide up the network and prevent the disease from spreading too far, no matter where it starts.

b) If you are allowed to remove just one node, which one should you remove to minimize the expected size of an outbreak in this network:



c) If you are allowed to remove two nodes, which should they be? How about three?

CHAPTER 16

EPIDEMICS ON NETWORKS

A discussion of the epidemic processes by which diseases spread over networks of contact between humans and animals

ONE of the reasons for the large investment the scientific community has made in the study of social networks is their connection to the spread of disease. Diseases spread over networks of contacts between individuals: airborne diseases like influenza and tuberculosis are communicated when two people breathe the air in the same room; other diseases and some parasites are communicated when people touch; sexually transmitted diseases such as HIV are communicated when people have sex. The patterns of such contacts can be represented as networks and a good deal of effort has been devoted to empirical studies of these networks' structure. We have already discussed some network aspects of epidemiology in the previous chapter when we considered site percolation as a model for the effects of vaccination. In this chapter we look in more detail at the connections between network structure and disease dynamics and at mathematical techniques that allow us to understand and predict the outcomes of epidemics.

There are also other network processes that can be viewed as generalized "spreading" processes akin to the spread of disease. The spread of news, rumors, or gossip through a population, for instance, has features in common with disease spread. Once a person hears a piece of news or information from an acquaintance they become capable of spreading that news to others—they have been "infected." The adoption of fashions or behaviors may also be

“contagious” in a similar way. The spread of information is arguably more complicated than the spread of a disease, in that information can also be spread by other means, including the Internet, the media, books, and so forth. Nonetheless, ideas and models for the spread of disease can be usefully applied to help us understand the spread of information.

Another variant of contagious behavior is *cascading failure*. Consider for instance the electrical power grid. Transmission lines regularly fail in the power grid, sometimes because of overload: if too much power is sent along a particular edge in the network then it can break down. When an edge fails, the power it previously carried must be rerouted some other way around the network in order to reach the consumers who need it, and in doing this we increase the load on other edges. This in turn can push those edges over the limit and cause them to fail too. The end result is a cascading failure in which a problem or fault in one initial spot in the network quickly spreads to a whole region. Cascading failures of this kind are the primary cause of large power outages and blackouts.¹

In this chapter we describe a range of methods and models for understanding and predicting the behavior of contagion processes and their generalizations on networks of various kinds.

16.1 MODELS OF THE SPREAD OF INFECTION

The biology of what happens when an individual (a “host” in the language of epidemiology) catches an infection is complicated. The pathogen responsible for the infection typically multiplies in the body while the immune system attempts to beat it back, often causing symptoms in the process. One or the other usually wins in the end, though sometimes neither, with the final result being the individual’s recovery, their death, or a chronic disease state of permanent infection. In theory if we want to understand fully how diseases spread through populations we need to take all of this biology into account, but in practice that is usually a dauntingly large job and it is rarely, if ever, attempted. Luckily there are more tractable approaches based on simplified models of disease spread that give a good guide to disease behavior in many cases and it is on these that we focus in this chapter.

¹It’s worth noting, however, that the causes and mechanisms of real-world power failures are usually quite complex, involving not just basic electrical principles, but also the combined effects of the actions of human operators and sophisticated control software whose behavior is not always ideal or easy to predict.

16.1.1 THE SI MODEL

In the typical mathematical representation of an epidemic the within-host dynamics of the disease is reduced to changes between a few basic disease states. Such states are sometimes called *compartments*, and mathematical models incorporating them are called *compartamental models*. In the simplest case there are just two states, *susceptible* and *infected*. An individual in the susceptible state is someone who does not have the disease yet but could catch it if they come into contact with someone who does. An individual in the infected state is someone who has the disease and can, potentially, pass it on if they come into contact with a susceptible individual.² Although this two-state classification sweeps a lot of biological details under the rug, it can capture broader features of disease dynamics and is a useful simplification in the present situation where we are focused more on what's happening at the level of networks and populations than on what's happening within the bodies of the individual population members.

Mathematical modeling of epidemics predates the study of networks by many years, stretching back at least as far as the pioneering work of Anderson McKendrick, a doctor and amateur mathematician who made foundational contributions to the field early in the twentieth century. The theories that he and others developed form the core of traditional mathematical epidemiology, which is an extensive and heavily researched field. Classic introductions to the subject include the highly theoretical 1975 book by Bailey [36] and the more recent and practically oriented book by Anderson and May [21]. The review article by Hethcote is also a good resource [232].

The traditional approach avoids discussing contact networks at all by making use of a *fully mixed* or *mass-action approximation*, in which it is assumed that every individual has an equal chance, per unit time, of coming into contact with every other—people mingle and meet completely at random in this approximation. This is, of course, not a realistic representation of the way the world is. In the real world, people have contact with only a small fraction of the population, and that fraction is not chosen at random, which is precisely

²In the epidemiology literature you may see the infected state referred to variously as infected, infectious, or infective. In most cases there is no difference between these terms; they are synonymous. One must be careful, however. As discussed later in the chapter, more sophisticated models of disease distinguish between a state in which an individual is infected with a disease but the infection has not yet developed to the point where the individual can pass it on, and a state where they can pass it on. The latter we called “infectious” or “infective”; the former we might call “infected,” though to avoid confusion it is more often called the “exposed” state. In the present simple two-state model, however, there is no difference between infected, infectious, and infective; they all describe the same state of having the disease and being able to pass it on.

why networks play an important role in the spread of disease and many other things. Nonetheless, a familiarity with the traditional approaches will be useful to us in our study of network epidemiology, so we will spend a little time looking at the basic principles.

Consider a disease spreading through a population of individuals. Let $S(t)$ be the number of individuals who are susceptible at time t and let $X(t)$ be the number who are infected.³ Technically, since the disease-spreading process is a random one, these numbers are not uniquely determined—if the disease were to spread through the same population more than once, even under very similar conditions, the numbers would probably be different each time. To get around this problem let us define S and X more carefully to be the average or expected numbers of susceptible and infected individuals,⁴ i.e., the numbers we would get if we ran the process many times under identical conditions and then averaged the results. Note that these average numbers will not, in general, be integers, even though the actual numbers of susceptible and infected individuals are necessarily always integers.

The number of infected individuals goes up when susceptible individuals contract the disease from infected ones. Suppose that people meet and make contact sufficient to result in the transmission of disease entirely at random with a per-individual rate β , meaning that each individual has, on average, β contacts with randomly chosen others per unit time.

The disease is transmitted only when an infected person has contact with a susceptible one. If the total population consists of n people, then the average probability that a person you meet at random is susceptible is S/n , and hence an infected person has contact with an average of $\beta S/n$ susceptible people per unit time. Since there are on average X infected individuals in total this means the overall average rate of new infections will be $\beta SX/n$ and we can write a differential equation for the rate of change of X thus:

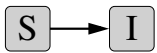
$$\frac{dX}{dt} = \beta \frac{SX}{n}. \tag{16.1}$$

At the same time the number of susceptible individuals goes down at the same rate:

$$\frac{dS}{dt} = -\beta \frac{SX}{n}. \tag{16.2}$$

³It might be more logical to use $I(t)$ for the number infected, and many authors do so, but we use X instead to avoid later confusion with the index i used to label nodes.

⁴For convenience we will usually drop the explicit t -dependence of $S(t)$ and $X(t)$ and, as here, just write S and X .



The allowed transitions between states can be represented by flow charts like this simple one for the SI model.

This simple mathematical model for the spread of a disease is called the fully mixed *susceptible–infected model*, or *SI model* for short.

It is often convenient to define variables representing the fraction of individuals who are in the susceptible and infected states thus:

$$s = \frac{S}{n}, \quad x = \frac{X}{n}, \quad (16.3)$$

in terms of which Eqs. (16.1) and (16.2) can be rewritten as

$$\frac{ds}{dt} = -\beta s x, \quad (16.4a)$$

$$\frac{dx}{dt} = \beta s x. \quad (16.4b)$$

In fact, we don't really need both of these equations, since it is also true that $S + X = n$ or equivalently $s + x = 1$ because every individual must be either susceptible or infected. With this condition it is easy to show that Eqs. (16.1) and (16.2) are really the same equation. Alternatively, we can eliminate s from (16.4) altogether by writing $s = 1 - x$, which gives

$$\frac{dx}{dt} = \beta(1 - x)x. \quad (16.5)$$

This equation, which occurs in many places in biology, physics, and elsewhere, is called the *logistic equation*. It can be solved using standard methods to give

$$x(t) = \frac{x_0 e^{\beta t}}{1 - x_0 + x_0 e^{\beta t}} \quad (16.6)$$

where x_0 is the value of x at $t = 0$. Generically, this produces an S-shaped “logistic growth curve” for the fraction of infected individuals, as shown in Fig. 16.1. The curve increases exponentially at short times, corresponding to an initial phase of the disease in which most of the population is susceptible, and then saturates as the number of susceptibles dwindles and the disease has a harder and harder time finding new victims.⁵

⁵There aren't many diseases that really saturate their population like this. Most real diseases that don't kill their victims are eventually defeated by the immune system. In addition, for many diseases some fraction of the population has a natural immunity that prevents them from being infected (meaning that when exposed to the pathogen their immune system sees it off so quickly that they never become infectious). And some diseases spread so slowly that a large fraction of the population never catches them because they die of other causes first. None of these phenomena is represented in this model.

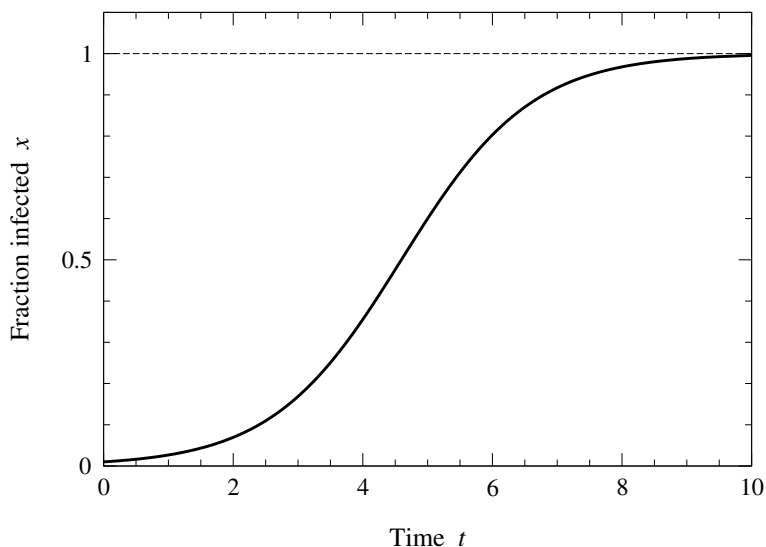


Figure 16.1: The classic logistic growth curve of the SI epidemic model. Starting from a small initial value (1% in this example) the number of infected individuals in an SI model grows exponentially at first, but growth eventually saturates as the supply of susceptible individuals is exhausted, and the curve levels off at $x = 1$ (dashed line).

16.1.2 THE SIR MODEL

The SI model is the simplest possible model of infection. There are many ways in which it can be extended to make it more realistic or more appropriate as a model of a specific disease. One common extension deals with recovery from disease.

In the SI model individuals, once infected, are infected (and infectious) forever. For many real diseases, however, people recover from infection after a certain amount of time because their immune system fights off the disease-causing agent. Furthermore, people often retain their immunity to the disease after such a recovery so that they cannot catch it again. To represent this behavior in our model we need a new third disease state, usually denoted R for *recovered*. The corresponding three-state model is called the *susceptible–infected–recovered* or *SIR model*.

With some diseases people do not recover but instead die. Although this is the complete opposite of recovery in human terms, it is essentially the same thing in epidemiological terms: it makes little difference to the disease whether a person is immune or dead—either way they are effectively removed from

the pool of potential hosts for the disease.⁶ Both recovery and death can be represented by the R state in our model. Diseases with mixed outcomes where people sometimes recover and sometimes die can also be modeled in this way—from a mathematical point of view we don't care whether the individuals in the R state are recovered or dead. For this reason some people say that the R stands for *removed* rather than recovered, so as to encompass both possibilities, and they refer to the corresponding model as the *susceptible–infected–removed* model.

The dynamics of the fully mixed SIR model has two stages. In the first stage, susceptible individuals become infected when they have contact with infected individuals. Contacts between individuals are assumed to happen at an average rate β per person as before. In the second stage, infected individuals recover (or die) at some constant average rate γ .

Given the value of γ we can calculate the length of time τ that an infected individual is likely to remain infected before they recover. The probability of recovering in any time interval $\delta\tau$ is $\gamma \delta\tau$ and the probability of not doing so is $1 - \gamma \delta\tau$. Thus the probability that the individual is still infected after a total time τ is given by

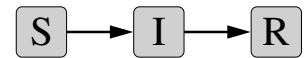
$$\lim_{\delta\tau \rightarrow 0} (1 - \gamma \delta\tau)^{\tau/\delta\tau} = e^{-\gamma\tau}, \tag{16.7}$$

and the probability $p(\tau) d\tau$ that the individual remains infected for time τ and then recovers in the interval between τ and $\tau + d\tau$ is this quantity times $\gamma d\tau$:

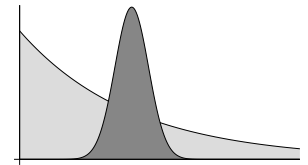
$$p(\tau) d\tau = \gamma e^{-\gamma\tau} d\tau, \tag{16.8}$$

which is just an exponential distribution. Thus, an infected person is most likely to recover immediately after becoming infected, but might in theory remain in the infected state for quite a long time—many times the mean infectious time (which is $1/\gamma$).

Neither of these behaviors is very realistic for most real diseases. With real diseases, most victims remain infected for about the same length of time, such as a week, say, or a month. Few stay in the infected state for much longer or shorter than the average (see figure). Nonetheless, we will for the moment stick



The flow chart for the SIR model.



For real diseases the distribution of times for which an individual remains infected is typically narrowly peaked around some average value, quite unlike the exponential distribution assumed by the SIR model.

⁶This is only approximately true. If people really do have a certain average number of contacts per unit time and assuming those contacts are with living people, then the presence of living but recovered people in the population reduces the number of contacts between infected and susceptible individuals. If, on the other hand, people die rather than recover from the disease then only susceptible and infected individuals are alive and the number of contacts between them will be correspondingly greater. In other words, as the population dwindles because of death from disease, the probability of contact between any two remaining people goes up. This effect can easily be incorporated into the model, but we don't do so here.

with this model because it makes the mathematics simple. This is one thing that will improve when we come to look at network models of epidemics.

Given these assumptions, the fractions s , x , and r of individuals in the three states are governed by the equations

$$\frac{ds}{dt} = -\beta sx, \tag{16.9a}$$

$$\frac{dx}{dt} = \beta sx - \gamma x, \tag{16.9b}$$

$$\frac{dr}{dt} = \gamma x, \tag{16.9c}$$

and in addition the three variables must satisfy

$$s + x + r = 1. \tag{16.10}$$

16.1.3 SOLUTION OF THE SIR MODEL

The SIR equations (16.9) can be solved as follows. First, we eliminate x between (16.9a) and (16.9c), giving

$$\frac{1}{s} \frac{ds}{dt} = -\frac{\beta}{\gamma} \frac{dr}{dt}. \tag{16.11}$$

Then we integrate both sides with respect to t to get

$$s = s_0 e^{-\beta r/\gamma}, \tag{16.12}$$

where s_0 is the value of s at $t = 0$ and we have chosen the constant of integration so that there are no individuals in the recovered state at $t = 0$. (Other choices are possible but we'll use this one for now.)

Now we put $x = 1 - s - r$ in Eq. (16.9c) and use Eq. (16.12) to get

$$\frac{dr}{dt} = \gamma(1 - r - s_0 e^{-\beta r/\gamma}). \tag{16.13}$$

If we can solve this equation for r , then we can find s from Eq. (16.12) and x from Eq. (16.10).

The solution is easy to write down in principle. It is given by

$$t = \frac{1}{\gamma} \int_0^r \frac{du}{1 - u - s_0 e^{-\beta u/\gamma}}. \tag{16.14}$$

Unfortunately, in practice we can't evaluate the integral in closed form. We can, however, evaluate it numerically. An example is shown in Fig. 16.2.

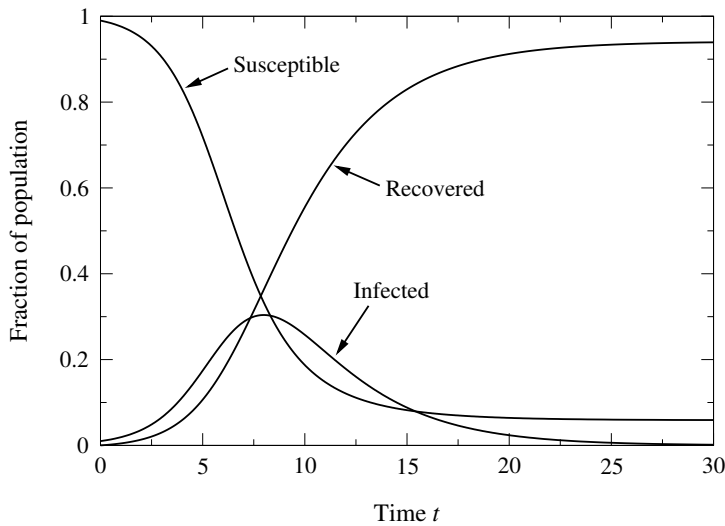


Figure 16.2: Time evolution of the SIR model. A numerical solution of the SIR equations (16.9). The three curves show the fractions of the population in the susceptible, infected, and recovered states as a function of time. The parameters are $\beta = 1$, $\gamma = 0.4$, $s_0 = 0.99$, $x_0 = 0.01$, and $r_0 = 0$.

There are a number of notable things about this figure. The fraction of susceptibles in the population decreases monotonically as susceptibles are infected and the fraction of recovered individuals increases monotonically. The fraction infected, however, goes up at first as people get infected, then down again as they recover, and eventually goes to zero as $t \rightarrow \infty$.

Note, however, that the number of susceptibles does not go to zero; the curve for $s(t)$ ends a little above the axis. This is because when $x \rightarrow 0$ there are no infected individuals left to infect the remaining susceptibles. Any individuals who survive to late enough times without being infected will probably never get the disease at all. They are the lucky ones who made it through the outbreak and out the other side. Similarly the fraction of recovered individuals does not quite reach one as $t \rightarrow \infty$.

The asymptotic value of r has an important practical interpretation: it is the total number of individuals who ever catch the disease during the entire course of the epidemic—the total size of the outbreak. It can be calculated from Eq. (16.13) as the value at which $dr/dt = 0$, which gives $r = 1 - s_0 e^{-\beta r/\gamma}$.

The initial conditions for the model can be chosen in a variety of ways, but the most common is to assume that the disease starts with either a single

infected individual or a small number c of individuals and everyone else in the susceptible state. In other words, the initial values of the variables are $s_0 = 1 - c/n$, $x_0 = c/n$, and $r_0 = 0$. In the limit of large population size $n \rightarrow \infty$, we can then write $s_0 \simeq 1$, and our final value of r satisfies

$$r = 1 - e^{-\beta r/\gamma}. \quad (16.15)$$

Interestingly, this is the same as the equation we derived in Section 11.5 for the size S of the giant component of a Poisson random graph, Eq. (11.16), provided we equate β/γ with the mean degree of the random graph, and this correspondence allows us immediately to say several useful things. First, we know what the size of the epidemic must look like as a function of the parameters β and γ : it will look like the plot of giant component size shown in Fig. 11.2b on page 351, except with β/γ along the horizontal axis instead of c . Second, it tells us that the size of the epidemic goes continuously to zero as β/γ approaches one from above. And for $\beta/\gamma \leq 1$, or equivalently $\beta \leq \gamma$, there is no epidemic at all. The simple explanation for this result is that if $\beta \leq \gamma$ then infected individuals recover faster than susceptible individuals become infected, so the disease cannot get a toehold in the population. The number of infected individuals, which starts small, goes down, not up, and the disease dies out instead of spreading.

The transition between the regimes where there is and is not an epidemic happens at the point $\beta = \gamma$, which is called the *epidemic threshold*. Note that there was no epidemic threshold in the simpler SI model: in that model the disease always spreads because individuals once infected never recover and hence the number of infected individuals cannot decrease.

One can think of the SI model as the special case of the SIR model in which $\gamma = 0$, so that β can never be less than γ .

16.1.4 BASIC REPRODUCTION NUMBER

An important quantity in the study of epidemics is the *basic reproduction number*, denoted R_0 , which is defined as follows. Consider the spread of a disease when it is just starting out, when there are only a few cases of the disease and the rest of the population is susceptible—what is called a *naive population* in epidemiology—and consider a susceptible person who catches the disease in this early stage of the outbreak. The basic reproduction number is defined to be the average number of additional people that such a person passes the disease on to before they recover. For instance, if each person catching the disease passes it on to two others on average, then $R_0 = 2$. If half of them pass it on to just one person and the rest to none at all, then $R_0 = \frac{1}{2}$, and so forth.

If we had $R_0 = 2$, then each person catching the disease would pass it on to two others on average, each of them would pass it on to two more, and

so forth, so that the number of new cases of the disease would double at each round, growing exponentially. Conversely if $R_0 = \frac{1}{2}$ the disease would die out exponentially. The point $R_0 = 1$ separates the regimes of growing and shrinking behavior and thus marks the epidemic threshold between cases where the disease multiplies and where it dies out.

We can calculate R_0 straightforwardly for the SIR model. If an individual remains infectious for time τ , then the expected number of others they will have contact with during that time is $\beta\tau$. The definition of R_0 is specifically for a naive population, and in a naive population all of the people with whom one has contact will be susceptible, and hence $\beta\tau$ is also the total number of people our infected individual will infect. Then we average over the distribution of τ , Eq. (16.8), to get the average number R_0 :

$$R_0 = \beta\gamma \int_0^\infty \tau e^{-\gamma\tau} d\tau = \frac{\beta}{\gamma}. \tag{16.16}$$

This gives us an alternative way of deriving the epidemic threshold for the SIR model: the threshold falls at $R_0 = 1$, i.e., at $\beta = \gamma$, the same result as we found in Section 16.1.3 by considering the long-time behavior.⁷

16.1.5 THE SIS MODEL

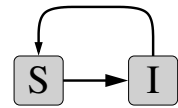
A different extension of the SI model is one that allows for *reinfection*, i.e., for diseases that don't confer immunity on their victims after recovery, or confer only limited immunity, so that individuals can be infected more than once. The simplest such model is the *SIS model*, in which there are just two states, susceptible and infected, and infected individuals move back into the susceptible state upon recovery. The differential equations for this model are

$$\frac{ds}{dt} = \gamma x - \beta s x, \tag{16.17a}$$

$$\frac{dx}{dt} = \beta s x - \gamma x, \tag{16.17b}$$

with

$$s + x = 1. \tag{16.18}$$



Flow chart for the SIS model.

⁷Note that when $\gamma = 0$, as in the SI model, Eq. (16.16) implies that $R_0 \rightarrow \infty$. This is because an infected individual remains infected indefinitely in the SI model and hence can infect an arbitrary number of others, so that R_0 is formally infinite. In any population of finite size, however, the number of people an individual infects will be finite.

Putting $s = 1 - x$ in Eq. (16.17b) gives

$$\frac{dx}{dt} = (\beta - \gamma - \beta x)x, \quad (16.19)$$

which has the solution

$$x(t) = (1 - \gamma/\beta) \frac{Ce^{(\beta-\gamma)t}}{1 + Ce^{(\beta-\gamma)t}}, \quad (16.20)$$

where the integration constant C is fixed by the initial value of x to be

$$C = \frac{\beta x_0}{\beta - \gamma - \beta x_0}. \quad (16.21)$$

In the case of a large population and a small number of initial carriers of the disease, x_0 is small and C is well approximated as $\beta x_0/(\beta - \gamma)$, which gives us

$$x(t) = x_0 \frac{(\beta - \gamma)e^{(\beta-\gamma)t}}{\beta - \gamma + \beta x_0 e^{(\beta-\gamma)t}}. \quad (16.22)$$

If $\beta > \gamma$ this produces a logistic growth curve similar to that of the basic SI model—see Fig. 16.3—but differing in one important respect: we never have the whole population infected with the disease. In the limit of long time the system finds a stable state where the rates at which individuals are infected and recover from infection are exactly equal and a steady fraction of the population—but not all of them—is always infected with the disease. (Which particular individuals are infected changes over time, however, as some recover and others are infected.) The fraction of infected individuals can be found from Eq. (16.22), or more directly from Eq. (16.19) by setting $dx/dt = 0$ to give $x = (\beta - \gamma)/\beta$. In the language of epidemiology the steady state is called an *endemic disease state*.

Note that the fraction infected in the endemic state goes to zero as β approaches γ , and if $\beta < \gamma$ then Eq. (16.22) predicts that the disease will die out exponentially. Thus, as in the SIR model, the point $\beta = \gamma$ marks an epidemic threshold between a state in which the disease spreads and one in which it dies out.⁸ As before, we can calculate a basic reproduction number R_0 , which again

⁸In a population of finite size there is always a non-zero probability that the disease will die out just by chance even in the regime above the epidemic threshold where the disease spreads. In a finite population, the fraction of individuals infected at any one time fluctuates because of randomness in the disease process, and if it ever fluctuates to zero—if by chance everyone with the disease recovers at the same time—then the disease no longer exists. The probability of this happening is extremely small for all but the tiniest of populations, but technically it is non-zero, so it's not strictly true to say simply that the disease always spreads when we are above the epidemic threshold. It will spread for a long time but eventually it will die out (by contrast with the situation below the threshold, where it dies out immediately).

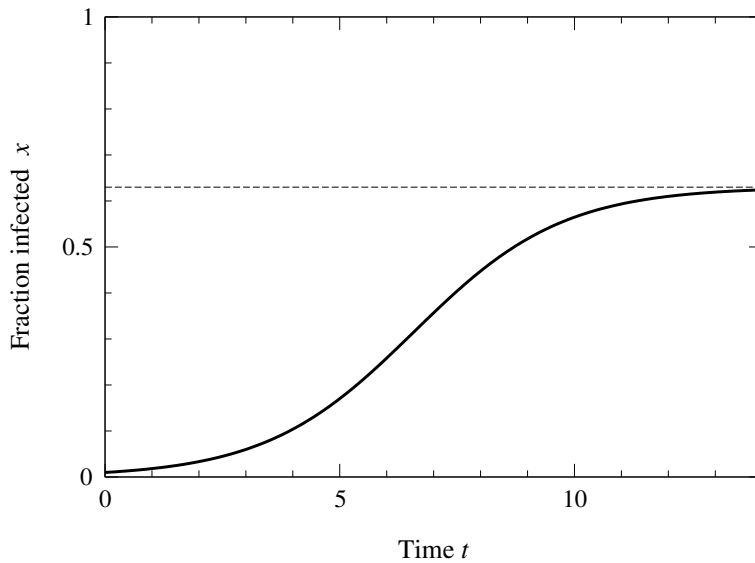
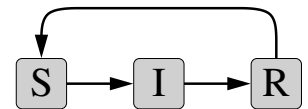


Figure 16.3: Fraction of infected individuals in the SIS model. The fraction of infected individuals in the SIS model grows with time following a logistic curve, as in the SI model. Unlike the SI model, however, the fraction infected never reaches unity, tending instead to an intermediate value (dashed line) at which the rates of infection and recovery are balanced. (Compare this figure with Fig. 16.1 for the SI model.)

takes the value $R_0 = \beta/\gamma$, giving us an alternative derivation of the position of the threshold as the point at which $R_0 = 1$.

16.1.6 THE SIRS MODEL

Another epidemic model, which we touch upon only briefly, is the *SIRS model*, another model incorporating reinfection. In this model individuals recover from infection and gain immunity as in the SIR model, but that immunity is only temporary. After a certain period of time individuals lose it and become susceptible again. We introduce a new parameter δ to represent the average



Flow chart for the SIRS model.

rate at which individuals lose immunity. Then the equations for the model are

$$\frac{ds}{dt} = \delta r - \beta sx, \quad (16.23a)$$

$$\frac{dx}{dt} = \beta sx - \gamma x, \quad (16.23b)$$

$$\frac{dr}{dt} = \gamma x - \delta r, \quad (16.23c)$$

and

$$s + x + r = 1. \quad (16.24)$$

The SIRS model cannot be solved analytically, although it can be treated using linear stability analysis and other tricks from the non-linear dynamics toolbox. One can also solve the differential equations numerically, which reveals that the SIRS model has a rich palette of behaviors depending on the values of the three parameters, including behaviors where the disease persists in an endemic state, where it dies out, and where it oscillates between outbreaks and periods of remission. We will not delve into the behavior of the SIRS model further in this chapter; the interested reader can find more details in Ref. [232].

16.1.7 OTHER EPIDEMIC MODELS

Many other epidemic models have been proposed to model the spread of particular types of diseases. Extra states can be introduced such as an “exposed” state that represents people who have caught a disease but whose infection has not yet developed to the point where they can pass it on to others. (The resulting model is sometimes called the SEIR model.) Another possibility is an initial immune state such that individuals start off immune to disease, lose that immunity and become susceptible, and from there move through the standard stages of infection. Models of this type, such as the MSIR model, are sometimes used to represent the maternally derived immunity that newborn babies possess (hence the letter M).

There are also models in which the population can grow as a result of births or immigration, models that distinguish between people who recover fully from disease and those who recover but remain carriers who can pass the disease to others, and many other variants. Comprehensive reviews of these and other models can be found in Refs. [21, 232].

16.1.8 COMBINATIONS OF DISEASES

All of the models we have described so far are models of the spread of a single disease. In the real world, however, there are many diseases circulating at the

same time, and sometimes many strains of the same disease. It is possible for diseases to interact in ways that change how they spread. Two cases in particular are common: cross-immunity and coinfection.

As we have discussed, many diseases impart upon their survivors immunity to further infection with the same disease. In some cases a disease can also impart immunity to *another* disease, a phenomenon known as *cross-immunity*. That is, after catching disease A a person is then immune to infection with some other disease B. A common case occurs when the diseases in question are both strains of the same larger disease family, such as two different varieties of the flu. Cross-immunity may be total (A imparts complete immunity to B) or partial (A reduces but does not eliminate the chance of infection with B). It may also be unilateral (A imparts immunity to B but not vice versa) or bilateral (either disease imparts immunity to the other).

Thus, for instance, one can make a version of the SI model that represents cross-immunity between two diseases A and B using four different disease states: susceptible (meaning a person has neither disease), infected with one disease or the other, and infected with both. The equations for the model would be

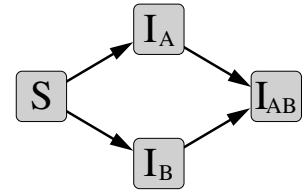
$$\frac{ds}{dt} = -\beta_{AS}(x_A + x_{AB}) - \beta_{BS}(x_B + x_{AB}), \quad (16.25a)$$

$$\frac{dx_A}{dt} = \beta_{AS}(x_A + x_{AB}) - \gamma_A x_A(x_B + x_{AB}), \quad (16.25b)$$

$$\frac{dx_B}{dt} = \beta_{BS}(x_B + x_{AB}) - \gamma_B x_B(x_A + x_{AB}), \quad (16.25c)$$

$$\frac{dx_{AB}}{dt} = \gamma_A x_A(x_B + x_{AB}) + \gamma_B x_B(x_A + x_{AB}). \quad (16.25d)$$

The parameters β_A and β_B measure the rate at which a susceptible person catches disease A or B respectively, while γ_A and γ_B measure the rate at which a person who currently has either A or B catches the other disease as well. In the simple case of complete cross-immunity we would have $\gamma_A = \gamma_B = 0$ and $x_{AB} = 0$ and the fourth equation would vanish, but in the more general case of partial immunity γ_A and γ_B would be non-zero and all four equations apply. Things become more complicated when we move to an SIR-type model that incorporates recovery from the diseases. In that case, as we have seen, there is an epidemic threshold that separates regimes in which a disease does or does not spread. Cross-immunity can then result in a disease that would otherwise spread falling below its epidemic threshold and dying out. If enough people become infected with disease A, for example, thereby giving those people immunity to disease B, then disease B may no longer be able to find enough susceptible hosts to spread.



Flow chart for the cross-immunity model.

Infection with one disease can also *enhance* the spread of another—the opposite of cross-immunity. That is, infection with disease A can make it easier, or in some cases possible at all, to get infected with disease B. As an example there are a number of infections that are found primarily in people with HIV: the immune deficiencies associated with HIV infection make possible other *coinfections* that otherwise would be unlikely to occur. An SI model mimicking this kind of behavior would have the same equations (16.25) as our model of cross-immunity, but the values of the parameters would be different. Instead of the β parameters being larger and the γ parameters being small (or zero), we would have β_A and γ_A large and β_B small (while γ_B can be either large or small).

When cross-immunity and coinfection processes are generalized to networks a range of interesting behaviors can happen, as we will see in Section 16.3.3.

16.1.9 COMPLEX CONTAGION AND THE SPREAD OF INFORMATION

As mentioned at the start of this chapter, the general ideas of disease spreading can be applied more broadly to shed light also on the spread of other things, such as information, news, rumors, or fashions. In some cases, the exact same models used for biological infections can be applied directly to other types of spreading process, but there are also some interesting variant models that capture unique features of the way information and other “social contagions” spread.

A classic example of a model of information spread is the *voter model*, originally proposed as a simple representation of the effects of peer pressure on voting behavior [303]. In this model it is assumed that a population of voters is asked to decide among some number of candidates for an election and at any particular moment every individual has a preferred candidate. One can think of these preferred candidates as being like the disease states in our earlier models. Initial preferences are typically assigned at random but individuals’ opinions can subsequently change following a simple algorithm: pairs of individuals meet at random and one of them adopts the preferred candidate of the other. Thus, in this model people’s opinions are affected only by other individuals. In real life many other factors play a role, such as news media, political campaigning, inherent preferences, and so forth, but the model ignores these.

The voter model is superficially similar to a kind of SI model: pairs of people meet and one of them “infects” the other with their opinion. On closer inspection, however, the behavior of the two models is very different. Consider the simplest case where voters are choosing between just two candidates. No

matter how many people currently prefer candidate 1, the probabilities of their number increasing or decreasing are identical: there is some probability that people with opposite opinions will meet and when they do one of them (chosen at random) infects the other, thereby either increasing or decreasing the ranks of candidate 1's supporters by one, with equal probability. Thus the model has no inherent preference for either candidate, even if one of them is currently more popular than the other, and the model's dynamics are therefore best thought of as a kind of random walk or random fluctuation that stops only when the model reaches *consensus*, meaning that everyone holds the same opinion and no further changes of opinion are possible.

In the form described here, therefore, the voter model is somewhat trivial. If one puts it on a network, however, so that individuals can only be "infected" by their neighbors and not by anyone, then it becomes more interesting [434]. The time to reach consensus depends strongly on the structure of the network. On networks with community structure, for instance, isolated groups of nodes can form with opinions contrary to the rest of the network, preventing the system from reaching consensus for long periods of time.

Another class of information spreading processes of particular interest are the *complex contagions*, those in which infection is communicated not by just a single contact between individuals but by multiple contacts or a particular pattern of contacts [98, 222, 464]. Imagine, for instance, that news of a particularly surprising or extraordinary event is spreading through a community. When you first hear the news, from a friend or co-worker perhaps, it seems incredible and you don't really believe it, or at least not enough to feel confident passing it along to anyone else. But then you hear it a second time from a different person and you realize that it's actually true, and at that point you start telling other people. In other words, you need to be *infected* twice before you become *infectious*.

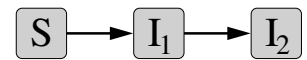
The equivalent of the SI model for this kind of complex contagion is a model in which there are three states: susceptible, meaning you have not heard the rumor, infected once, and infected twice (and hence infectious). The equations would be

$$\frac{ds}{dt} = -\beta s x_2, \quad (16.26a)$$

$$\frac{dx_1}{dt} = \beta s x_2 - \gamma x_1 x_2, \quad (16.26b)$$

$$\frac{dx_2}{dt} = \gamma x_1 x_2. \quad (16.26c)$$

Note how all terms involve x_2 —you can only catch the infection from individuals in the twice-infected state, since those in the once-infected state do not



Flow chart for a two-step complex contagion process.

spread it.

One could also create variants of this model in which an individual must be infected three or more times before they themselves become infectious, or ones where the required number of infections varies from one individual to another. The number could, for instance, be drawn at random for each individual from some specified distribution. One can also make SIR-style versions in which there is a recovered state, or many other variations on the basic theme.

16.2 EPIDEMIC MODELS ON NETWORKS

As discussed in Section 16.1.1, the traditional approach to epidemic modeling assumes “full mixing” of the population, meaning that any individual can have contact with and potentially transmit disease to any other. In the real world, however, this is not a realistic assumption. Most people have a set of regular acquaintances—friends and family, neighbors, co-workers, and so forth—whom they meet with some regularity and the rest of the population can safely be ignored. The pattern of people’s regular contacts can be represented as a network and the structure of that network can have a strong effect on the way a disease spreads.

Network models of disease work in basically the same way as the fully mixed models we have already seen but make use of this network of contacts instead of assuming that contact is possible with the entire population. Let us define the *transmission rate* or *infection rate* for our network disease process to be the probability per unit time that infection will be transmitted between two individuals, one susceptible and one infected, who are connected by an edge in the appropriate network. To put that another way, it is the rate at which contact sufficient to spread the disease occurs between any two individuals connected by an edge. The transmission rate is commonly denoted β by analogy with the corresponding quantity appearing in the fully mixed models (see Section 16.1.1), and we will adopt that notation here, although you should note that the two parameters are not exactly equivalent since β in the fully mixed case is the rate of contacts between an infected individual and all others in the population, whereas in the network case it is the rate of contact with just one other.

The transmission rate is partly a property of the disease being considered: some diseases are transmitted more easily than others and so have higher transmission rates. But transmission rate is also a property of the social and behavioral parameters of the population. In some communities or cultures, for instance, people may meet face-to-face with their acquaintances more often than in others. Some communities may favor more close physical contact be-

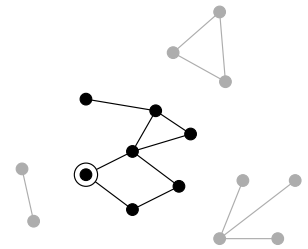
tween acquaintances than others. Such differences could produce a significant difference in transmission rates.

16.3 OUTBREAK SIZES AND PERCOLATION

Given a value for the transmission rate one can define models for the spread of disease over a network. Each of the models introduced in the first part of this chapter can be generalized to the network case. Consider the SI model, for instance. In the network version of this model we have n individuals represented by the nodes of our network, with most of them in the susceptible state at time $t = 0$ and just a small fraction x_0 , or maybe even just a single node, in the infected state. With probability β per unit time, infected nodes spread the disease to their susceptible neighbors, and infected nodes remain infected forever, so over time the disease spreads across the network.

It is difficult to solve a model such as this for a general network, and in many cases the best we can do is simulate it on a computer. There is, however, one important quantity that is straightforward to calculate, namely the total size of the disease outbreak. It is clear that in the limit of long time in this model every individual who *can* be infected by the disease *is* infected: since infected individuals remain infected forever, their susceptible neighbors will always, in the end, also become infected, no matter how small the transmission rate, so long as it is not zero. The only condition for a node to be infected therefore is that there must be at least one path to that node through the network from one of the initially infected individuals. This means that in the limit of long time the disease will spread from every initial carrier to infect all nodes in the component to which the carrier belongs. In the simplest case, where the disease starts out with a single infected carrier, just one component will be infected.

As we have seen, however, most networks have one large component that contains a significant fraction of all nodes in the network, plus, typically, a selection of smaller components. If we have this kind of structure then an interesting behavior emerges. If we start with a single infected individual, and if that individual turns out to belong to the large component, then the disease will infect the large component and we will have a large outbreak. If the individual belongs to one of the small components, however, the disease will only infect the few members of that small component and then die out. If the initial carrier of the disease is chosen uniformly at random from the network, then the probability that it will fall in the large component and we get a large outbreak is simply equal to the fraction S of the network occupied by the large component, and the size of the outbreak as a fraction of the network will also be S . Conversely, with probability $1 - S$ the initial carrier will fall in one of the



An outbreak starting with a single infected individual (circled) will eventually affect all those in the same component of the network, but leave other components untouched.

small components and the outbreak will be small, with size given by the size of the appropriate small component.

This constitutes a new type of behavior not seen in fully mixed models. In fully mixed models the possible behaviors are also either a run-away epidemic that affects a large fraction of the population, or an outbreak that affects only a few then dies out. But the choice between these outcomes was uniquely determined by the choice of model and the model parameters. For a given model and set of parameter values the disease always did either one thing or the other. In our network model, however, the behavior depends on the network structure and on the position in the network of the first infected individual. Thus there is a new stochastic element in the process: with identical model parameters and an identical network the disease sometimes takes off and sometimes dies out.

16.3.1 OUTBREAK SIZES IN THE SIR MODEL

The situation becomes more interesting still when we look at the SIR model. In the SIR model individuals remain infectious for only a finite amount of time and then they recover, so it is in general no longer true (as in the SI model) that the susceptible neighbor of an infected individual will always get infected in the end. If they are lucky such neighbors may never catch the disease. The probability of this happening—the probability that the disease is not transmitted—can be calculated in a manner similar to the calculation of Eq. (16.7) and is equal to $e^{-\beta\tau}$, where β is again the transmission rate and τ is the amount of time for which the infected individual remains infected. Thus the probability that the disease *is* transmitted is

$$\phi = 1 - e^{-\beta\tau}. \quad (16.27)$$

We will refer to this quantity as the *transmission probability*.

For simplicity, let us suppose that every infected individual remains infectious for the same length of time. This differs from the fully mixed version of the model, where τ was distributed according to an exponential distribution (see Eq. (16.8)), but in many cases is actually more realistic. As mentioned in Section 16.1.2, observed values of τ for many diseases are narrowly concentrated about a mean value, and their distribution is far from being exponential.

With this assumption, the transmission probability ϕ is a constant across the network. Every susceptible individual has equal probability of catching the disease from their infected neighbor (though if they have more than one infected neighbor the total probability is higher). Note that if the transmission probability is 1, then we recover the behavior of the SI model—every node that

can be infected is infected. Thus we can consider the SI model to be the special case of the SIR model where $\phi = 1$.

Now here is a useful trick, originally introduced by Mollison [336] and Grassberger [215]. Let us take our network and “fill in” or “occupy” each edge with probability ϕ , or not with probability $1 - \phi$. This is exactly the bond percolation process introduced in Section 15.1, where a fraction ϕ of edges are occupied uniformly at random. The occupied edges represent those along which disease will be transmitted if it reaches either of the nodes at the ends of the edge. That is, the occupied edges represent contacts sufficient to spread the disease, but not necessarily actual disease transmission: if the disease doesn’t reach either end of an edge, then it will not be transmitted along that edge.

With this in mind consider now the spread of a disease that starts at a randomly chosen node. We can immediately see that the set of nodes to which the disease will ultimately spread is precisely the set reachable from the initial node via paths that go along occupied edges—the disease simply passes from one node to another by traversing occupied edges until all reachable nodes have been infected. The end result is that the disease infects all members of the bond percolation cluster to which the initial carrier belongs.

It is important to appreciate that, as with our treatment of the network version of the SI model in the previous section, this process does not give us any information about the temporal evolution of the disease outbreak. Individual infection events are stochastic and a calculation of the curve of infections as a function of time requires a more complicated analysis that takes their randomness into account. However, if we want to know only about long-time behavior, about the overall total number of individuals infected by the disease, then all we need do is count the nodes in the appropriate percolation cluster.

Bond percolation is in many ways similar to the site percolation processes we studied in Chapter 15. Consider Fig. 16.4. For low edge occupation probability (Fig. 16.4a) there are just a few occupied bonds which form small clusters disconnected from one another. But as ϕ increases there comes a point, the percolation transition, where the disconnected clusters grow large enough to join together and form a giant cluster, although usually there are also other small clusters that are not joined to the giant cluster (Fig. 16.4b). As ϕ increases still further, the giant cluster grows, reaching its maximum size when $\phi = 1$ (Fig. 16.4c). Note, however, that this maximum size is not generally equal to the size of the whole network. Even when every edge in the network is occupied, the size of the largest cluster is still limited to the size of the largest component on the network, which is usually smaller than the whole network.

Translating these ideas into the language of epidemiology, we see that for small values of ϕ the cluster to which the initial carrier of a disease belongs

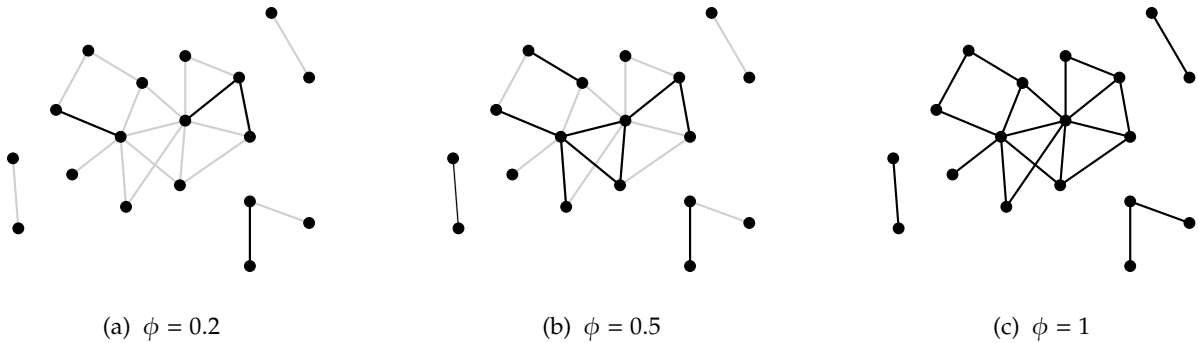


Figure 16.4: Bond percolation. In bond percolation, a fraction ϕ of the edges in a network are filled in or “occupied” at random to create connected clusters of nodes. (a) For small occupation probability ϕ the clusters are small. (b) Above the percolation threshold a large cluster forms, though there are usually still some small clusters as well. (c) When $\phi = 1$ all edges are occupied but the large cluster may still not fill the whole network: at $\phi = 1$ the largest cluster corresponds to the largest component of the network, which is often just a subset of the whole network.

must be small, since all clusters are small. Thus, in this regime we will have only a small disease outbreak and most members of the population will be uninfected. Once we reach the percolation transition, however, and a giant cluster forms, then a large outbreak of the disease—an epidemic—becomes possible, although not guaranteed. If the giant cluster of the percolation process occupies a fraction S of the entire network, then our randomly chosen initial carrier will fall within it with probability S , and if it does then the disease will spread to infect the whole giant cluster, creating an epidemic reaching a fraction of the population also equal to S . With probability $1 - S$, on the other hand, the initial carrier will fall in one of the small clusters and we will have only a small outbreak of the disease. As ϕ increases, S also increases and hence both the probability and the size of an epidemic increase with ϕ .

Thus, the percolation threshold for bond percolation on our network corresponds precisely to the epidemic threshold for an SIR-type disease on the same network, where the edge occupation probability is equal to the transmission probability ϕ of the disease, which is given in terms of the transmission rate β and recovery time τ by Eq. (16.27), and the sizes of outbreaks are given by the sizes of the bond percolation clusters. This mapping between percolation and epidemics is a powerful one that gives us a way to understand the effects of network structure on the spread of disease.

It is important to note that even when the value of ϕ is above the epidemic threshold we are not guaranteed that there will be an epidemic. This is similar

to the situation we saw for the simpler SI model on a network, but different from the situation for the fully mixed (non-network) SIR model of Section 16.1.2, where an epidemic always takes place if we are above the epidemic threshold. In many ways the behavior of our network model is more realistic than that of the fully mixed model. In real life it is true that outbreaks do not always result in epidemics. A disease can die out because, just by chance, its first victims happen not to pass the disease on to others. Our theory tells us that the probability of this happening is $1 - S$, where S is the size of the giant cluster for bond percolation on the network, which is also the size of the epidemic if it does happen. The value of $1 - S$ is usually small when we are well above the epidemic threshold, but can be quite large if we are only a little above the threshold, meaning that the probability of the disease dying out can be large in this regime.

It is also important to bear in mind that percolation is a stochastic process. We occupy edges at random on our network to represent the random nature of the contacts that transmit the disease. Two outbreaks happening under the same conditions on the same network would not necessarily spread along the same edges and the shapes of the percolation clusters would not necessarily be the same. Thus, a node that happens to belong to the giant cluster on one occasion might not belong to it on another and our theory cannot make exact predictions about what will happen when. The best we can do is calculate probabilities or average behaviors. We can, for instance, calculate the expected number of people who would be affected by an outbreak, but we cannot predict the exact number for any given outbreak.

16.3.2 SIR MODEL AND THE CONFIGURATION MODEL

In Section 15.2.1 we showed that it is possible to calculate exactly the average behavior of a site percolation process on configuration model networks. With only slight modifications the same approach can also be used for bond percolation and hence we can calculate the size distribution of epidemics and the position of the epidemic threshold in such networks.

Consider an SIR epidemic of the kind discussed in the previous section, taking place on a configuration model network with degree distribution p_k , and consider the corresponding bond percolation process with edge occupation probability ϕ given by Eq. (16.27). Let u be the average probability that a node is not connected to the giant percolation cluster via a specific one of its edges. There are two ways this can happen: either the edge in question is unoccupied (with probability $1 - \phi$), or it is occupied (probability ϕ) but the node at the other end of the edge is itself not a member of the giant cluster. The latter

happens only if that node is not connected to the giant cluster via any of its other edges, which happens with probability u^k if there are k such edges. Thus the total probability is $1 - \phi + \phi u^k$.

The value of k is distributed according to the excess degree distribution

$$q_k = \frac{(k+1)p_{k+1}}{\langle k \rangle} \quad (16.28)$$

(see Eq. (12.16)) and, averaging over k , we then arrive at a self-consistent expression for u thus:

$$u = \sum_{k=0}^{\infty} q_k (1 - \phi + \phi u^k) = 1 - \phi + \phi \sum_{k=0}^{\infty} q_k u^k = 1 - \phi + \phi g_1(u), \quad (16.29)$$

where g_1 is the probability generating function for the excess degree distribution, defined in Eq. (12.97). Equation (16.29) is the same as the corresponding equation for the site percolation case, Eq. (15.4), and has the same solutions.

The probability that a node of total degree k does not belong to the giant cluster is now simply u^k , and the average such probability over the whole network, which is equal to $1 - S$, is calculated by averaging u^k over the degree distribution p_k , giving

$$S = 1 - \sum_{k=0}^{\infty} p_k u^k = 1 - g_0(u), \quad (16.30)$$

where g_0 is the generating function for the degree distribution, Eq. (12.96). Equation (16.30) differs from the corresponding equation for the site percolation case, Eq. (15.2), by an overall factor of ϕ , but is otherwise the same. Thus the shape of the curve for S as a function of ϕ will be different from the site percolation case, but the position ϕ_c of the percolation transition, which is dictated by the solution of Eq. (16.29), will be the same. The solution of Eq. (16.29) was shown graphically in Fig. 15.2 and the position of the transition is given by Eq. (15.7) to be

$$\phi_c = \frac{1}{g_1'(1)} = \frac{\langle k \rangle}{\langle k^2 \rangle - \langle k \rangle}. \quad (16.31)$$

This equation thus also gives us the position of the epidemic threshold in terms of the transmission probability ϕ . If we prefer our solution in terms of the more fundamental parameters β and τ , we can rearrange Eq. (16.27) to give

$$\beta\tau = -\ln(1 - \phi_c) = \ln \frac{\langle k^2 \rangle - \langle k \rangle}{\langle k^2 \rangle - 2\langle k \rangle}. \quad (16.32)$$

If $\beta\tau$ exceeds this value then there is the possibility of an epidemic, though not the certainty, since the initial carrier or carriers of the disease could by chance fall outside the giant cluster. If $\beta\tau$ is smaller than this value then an epidemic is impossible, no matter where the initial carrier is. The probability of the epidemic, if one is possible, is given by S , Eq. (16.30), as is the size of the epidemic if and when one occurs.

Since the epidemic behavior of the model is controlled by the combination of parameters $\beta\tau$, the epidemic transition can be driven either by an increase in the infectiousness time τ , which is a property of the particular disease under study, or by an increase in the transmission rate β , which is a property both of the disease and of the behavior of members of the population. At the same time, the precise position of the transition in terms of these variables, as well as the probability and size of any epidemic that occurs, depends on the structure of the network via the moments $\langle k \rangle$ and $\langle k^2 \rangle$ of the degree distribution. This contrasts with the fully mixed model of Section 16.1.2, which incorporated no network effects.

Because of the close similarity between the site and bond percolation problems, we can easily translate a number of the results of Chapter 15 into the language of epidemics. For instance, a random graph with a Poisson degree distribution with mean c , which has $g_0(z) = g_1(z) = e^{c(z-1)}$, has an epidemic threshold that falls at $\phi_c = 1/c$ (Eq. (15.11)), or

$$\beta\tau = \ln \frac{c}{c-1}, \quad (16.33)$$

and the size of the epidemic, when there is one, is given by the solution to the equations

$$u = 1 - \phi + \phi e^{c(u-1)}, \quad (16.34)$$

$$S = 1 - e^{c(u-1)}. \quad (16.35)$$

The first of these equations can be rearranged to read $1 - u = \phi(1 - e^{c(u-1)}) = \phi S$ and substituting into the second then gives

$$S = 1 - e^{-\phi c S}, \quad (16.36)$$

which has no simple closed-form solution,⁹ but can easily be solved numerically by making an initial guess at the solution ($S = \frac{1}{2}$ seems to work well) and then iterating the equation to convergence. Figure 16.5 shows the results for the

⁹The solution can be written in closed form using the *Lambert W-function*, which is defined to be the solution of the equation $W(z)e^{W(z)} = z$. In terms of this function, the size of the epidemic is

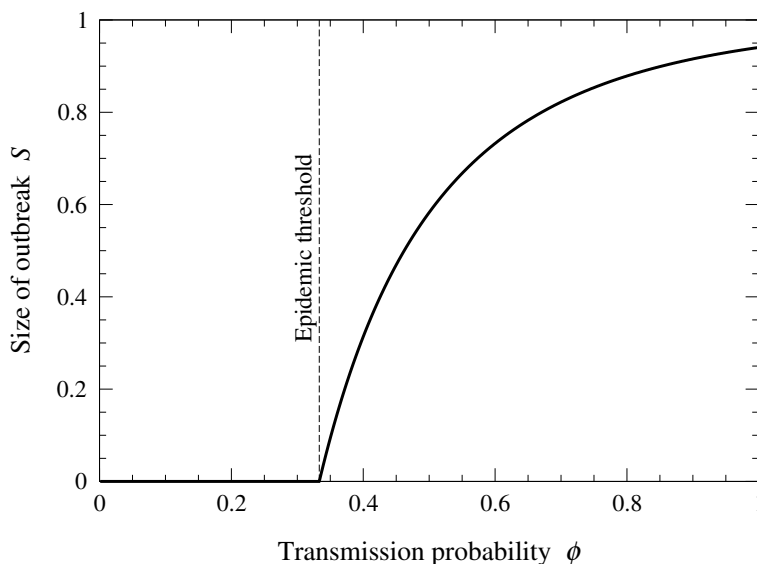


Figure 16.5: Size of an epidemic on a Poisson random graph. The size of an epidemic outbreak of an SIR-type disease on a Poisson random graph with mean degree $c = 3$, calculated by numerical solution of Eq. (16.36). The vertical dashed line marks the position of the epidemic threshold at $\phi_c = 1/c$.

case $c = 3$. The curve shows the size of the epidemic as a function of ϕ , starting out at zero on the left-hand side of the graph where ϕ is too low to support spread of the disease, then growing continuously once we pass the epidemic threshold.

Note that (16.36) is similar in form to Eq. (16.15) for the fully mixed model, but with different parameters. The similarity is not coincidental. In the fully mixed model an infected individual infects others chosen uniformly at random from the population, and in the Poisson random graph the network neighbors of any individual are also chosen uniformly at random. It is possible to show

given by

$$S = 1 + \frac{W(-\phi c e^{-\phi c})}{\phi c}.$$

Alternatively, we can rearrange Eq. (16.36) to give ϕ as a function of S rather than the other way around:

$$\phi = -\frac{\ln(1-S)}{cS}.$$

This expression can be useful for making plots of S .

that there is a direct correspondence between the traditional fully mixed model and the network model on a random graph [44].¹⁰

Another important case is the scale-free network with its power-law degree distribution. As we saw in Section 15.2.1, if the exponent α of the power law in such a network lies in the usual range $2 < \alpha < 3$, then $\phi_c = 0$, because $\langle k^2 \rangle$ diverges while $\langle k \rangle$ remains constant and hence Eq. (16.31) goes to zero. Thus in the power-law case there is always an epidemic, no matter how small the probability of transmission of the disease, at least in the limit of infinite network size. (For finite networks, $\langle k^2 \rangle$ is not infinite, but very large, and ϕ_c is correspondingly very small, but not precisely zero.)

This statement is, however, slightly misleading since, as discussed at the end of Section 15.2.1, the size of the giant cluster in a scale-free network becomes very small as we approach $\phi = 0$; it generally decays faster than linearly with ϕ . Thus, although technically there may be an epidemic for all positive values of ϕ , it can be very small in practice, affecting only the tiniest fraction of the population. (On the other hand, the difference between non-epidemic behavior and epidemic behavior, even with a tiny value of S , will become crucial when we look at models such as the SIS model that incorporate reinfection. In such models the epidemic threshold separates the regime in which the disease persists and the regime in which it becomes extinct, an important distinction even if the number of individuals infected is small.)

16.3.3 COEXISTING DISEASES

The connection between the SIR model and percolation can be extended to shed light on other, more complicated disease processes as well. As an example, consider the case, discussed in Section 16.1.8, of two SIR-type diseases spreading through the same population, one of which confers on its survivors (cross-)immunity to the other. Suppose, for the sake of simplicity, that disease 1 passes completely through the population and dies out before disease 2 starts to spread. In that case, the spread of disease 1 is just a normal SIR process that can be mimicked using bond percolation as discussed in the previous sections. The net result of this process is that some fraction of the nodes in the network are infected with disease 1, then recover again and acquire immunity, not only

¹⁰The differences in parameters arise because we are considering a slightly different disease process (one in which each individual is infectious for the same amount of time, rather than the exponential distribution used in the fully mixed model), and also because in the network model β is the transmission rate per edge, rather than the rate for the whole network—this is what gives us the factor of c in the exponent of Eq. (16.36).

to disease 1 itself but also to disease 2. We will assume that the cross-immunity is total in this case—survivors of disease 1 are completely immune to disease 2.

Thus the passage of disease 1 through the network effectively vaccinates all nodes in the percolation cluster in which it starts against disease 2. As far as disease 2 is concerned, those nodes are now removed from consideration—see Fig. 16.6—and the second disease is left to spread on the remaining nodes. The spread of the second disease can now be modeled as another percolation process on the network formed by these remaining nodes and their connecting edges.

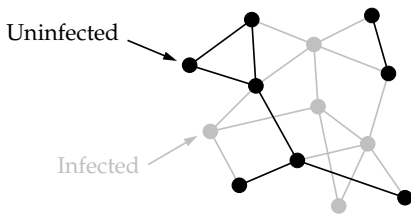


Figure 16.6: Cross-immunity on a network. When a disease imparts cross-immunity to the nodes it infects, those nodes effectively become removed from the network (gray nodes in this figure). A subsequent outbreak will then only spread if the remaining uninfected nodes form a giant percolating cluster.

Consider, for instance, the configuration model again. In general our two diseases can have different transmission probabilities, which we will denote ϕ_1 and ϕ_2 . Let us suppose that the transmission probability ϕ_1 for disease 1 is large enough to put it above the epidemic threshold, and moreover that an epidemic does actually occur, meaning that the initial carrier of the disease falls in the giant cluster. Then the first disease will effectively remove from the network all nodes in that giant cluster, immunizing them against the second disease. The remaining nodes in the network still form a configuration model network—the probability of any two of them being connected by an edge is exactly the same as it always was. However—and this is crucial—the degree distribution of this “residual network” will not be the same as the original network we started with, because nodes with higher degree are more

likely to be infected by disease 1, and hence removed from consideration. So the residual network will contain a preponderance of lower-degree nodes, a fact we must take into account if we are to accurately calculate the spread of disease 2. The calculation, which was first given in Ref. [358], goes as follows.

Let $P(\text{uninf}, m|k)$ be the probability that a node does not catch the first disease, so that it falls in the residual network, and that it has degree m within that residual network, given that it had total degree k to begin with. Then, by definition, the node has m edges (occupied or not) to other nodes that did not catch disease 1, plus $k - m$ edges to nodes that did, but the latter edges must be unoccupied (since if any of them were occupied then our node would have caught disease 1). In the notation of Section 16.3.2, the probability that the node at the other end of an edge is uninfected is u^k , or $\sum_k q_k u^k = g_1(u)$ after we average over the excess degree k . Meanwhile, the probability that an edge

is unoccupied and leads to a node that is infected with disease 1 is

$$(1 - \phi_1)[1 - g_1(u)] = 1 - \phi_1 - g_1(u) + \phi_1 g_1(u) = u - g_1(u), \quad (16.37)$$

where we have made use of Eq. (16.29). Putting these probabilities together, we have

$$P(\text{uninf}, m|k) = \binom{k}{m} [g_1(u)]^m [u - g_1(u)]^{k-m}, \quad (16.38)$$

where the factor of $\binom{k}{m}$ accounts for the different ways of choosing m edges from k possibilities.

Now we multiply this probability by the probability p_k that a node has degree k , and divide by the probability $P(\text{uninf})$ that the node is uninfected to begin with, which is just $1 - S = g_0(u)$ in the notation of Section 16.3.2. Then, applying Bayes' rule, we have

$$P(m, k|\text{uninf}) = \frac{p_k}{P(\text{uninf})} P(\text{uninf}, m|k) = \frac{p_k}{g_0(u)} \binom{k}{m} [g_1(u)]^m [u - g_1(u)]^{k-m}. \quad (16.39)$$

Summing this probability over all values of k from m upward (since $k \geq m$ necessarily), we now get the degree distribution of the residual network of uninfected nodes:

$$P(m|\text{uninf}) = \frac{1}{g_0(u)} \sum_{k=m}^{\infty} p_k \binom{k}{m} [g_1(u)]^m [u - g_1(u)]^{k-m}. \quad (16.40)$$

The generating function for this degree distribution is

$$\begin{aligned} f_0(z) &= \sum_{m=0}^{\infty} P(m|\text{uninf}) z^m = \frac{1}{g_0(u)} \sum_{m=0}^{\infty} \sum_{k=m}^{\infty} p_k \binom{k}{m} [g_1(u)]^m [u - g_1(u)]^{k-m} z^m \\ &= \frac{1}{g_0(u)} \sum_{k=0}^{\infty} p_k \sum_{m=0}^k \binom{k}{m} [g_1(u)]^m [u - g_1(u)]^{k-m} z^m \\ &= \frac{1}{g_0(u)} \sum_{k=0}^{\infty} p_k [u + (z-1)g_1(u)]^k = \frac{g_0(u + (z-1)g_1(u))}{g_0(u)}. \end{aligned} \quad (16.41)$$

Once we have this generating function we can also calculate the generating function $f_1(z)$ for the corresponding excess degree distribution from Eq. (12.33):

$$f_1(z) = \frac{f_0'(z)}{f_0'(1)} = \frac{g_1(u + (z-1)g_1(u))}{g_1(u)}. \quad (16.42)$$

Armed with these two generating functions we can now calculate a variety of features of the spread of the second disease. For instance, by analogy with

Eqs. (16.29) and (16.30), the fraction C of the residual network infected in an epidemic outbreak of the second disease is given by

$$C = 1 - f_0(v), \quad v = 1 - \phi_2 + \phi_2 f_1(v). \quad (16.43)$$

We can also calculate the position of the epidemic threshold for the second disease, which is given by the equivalent of Eq. (16.31):

$$\phi_2 = \frac{1}{f_1'(1)} = \frac{1}{g_1'(u)}, \quad (16.44)$$

where we have used Eq. (16.42) to derive the second equality.

But now we make an interesting observation. If the value of $g_1'(u)$ is less than 1, then Eq. (16.44) implies that ϕ_2 would have to be greater than 1 for disease 2 to spread, which is not possible since ϕ_2 is a probability. In this regime, therefore, the second disease will never spread. An equivalent way to say the same thing is that in this regime the first disease infects so many nodes that the residual network it leaves behind does not have a giant component at all, in which case the second disease obviously cannot spread. To see this, recall that in Section 12.6.2 we showed that a network has a giant component if $g_1'(1) > 1$ (see Eq. (12.42)). The equivalent statement for our residual network is that there is a giant component if $f_1'(1) > 1$, or equivalently $g_1'(u) > 1$, and hence it does not have a giant component if $g_1'(u) < 1$.¹¹

Now consider the behavior of our two diseases if we take a network and slowly increase the value of ϕ_1 from zero, while keeping the value of ϕ_2 fixed. Initially, the first disease will not spread, since its probability of transmission is too low. Provided ϕ_2 is large enough, however, disease 2 will spread. As ϕ_1 increases, it will eventually reach the epidemic threshold given by Eq. (16.31) and disease 1 will start to spread. When it does so, it will impart immunity on the nodes it infects and hence reduce the number of nodes available to be infected by disease 2. At this point we will start to see the size of the outbreak of disease 2 decreasing.

As we increase ϕ_1 still further the size of the outbreak of disease 2 will continue to decrease and eventually will reach zero at the point where disease 1 inoculates so many nodes that there is no longer a giant component among the nodes it leaves behind, and disease 2 can no longer spread. Thus there is a second threshold as a function of ϕ_1 in addition to the ordinary epidemic threshold, sometimes called the *coexistence threshold*. The coexistence threshold

¹¹The borderline case where $g_1'(u) = 1$ is a marginal situation where technically there is no giant component but the largest component has size going as a power of n .

is the point at which disease 1 infects so many nodes that disease 2 can no longer spread and hence the two diseases cannot both exist in the same population. We can calculate the position of the coexistence threshold by solving the condition $g'_1(u) = 1$ for the value of u and then rearranging Eq. (16.29) to give the corresponding value of ϕ_1 :

$$\phi_1 = \frac{1 - u}{1 - g_1(u)}. \quad (16.45)$$

For example, consider again the case of a network with a Poisson degree distribution with mean c , so that $g_0(z) = g_1(z) = e^{c(z-1)}$, as in Section 16.3.2. Then, as we showed there, the epidemic threshold for the first disease falls at $\phi_1 = 1/c$. The equation $g'_1(u) = 1$ takes the form $ce^{c(u-1)} = 1$, which can be written either as $1 - u = (1/c) \ln c$ or as $g_1(u) = 1/c$. Substituting these forms into Eq. (16.45) then gives the position of the coexistence threshold:

$$\phi_1 = \frac{\ln c}{c - 1}. \quad (16.46)$$

Figure 16.7 shows the sizes of the outbreaks of the two diseases for the case $c = 3$, as a function of ϕ_1 , with $\phi_2 = 1$. The behavior is exactly as we expect. From left to right in the figure we start in a regime where disease 1 does not spread but disease 2 does, pass into a coexistence region (shaded) where both diseases spread, and then into a region where only disease 1 spreads. Because we chose the maximum possible value of $\phi_2 = 1$, disease 2 persists right up to the coexistence threshold. If we had chosen a smaller value, disease 2 would have died out at an earlier point. Thus the coexistence threshold represents the point beyond which disease 2 cannot spread, no matter what its transmission probability is, and not necessarily the point at which it actually dies out.

16.3.4 COINFECTION

As described in Section 16.1.8, two diseases spreading in a single population need not necessarily hinder one another. In some cases one disease can facilitate the spread of another. We can also model this situation using the percolation techniques developed here, but the calculations are significantly more complex than for the cross-immunity case studied in the previous section. Consider, for instance, the simplest case, analogous to that of Section 16.3.3, in which two diseases spread through a population, one after another, and an individual can only catch the second disease if they have already been infected with the first. Thus the first disease effectively selects a subset of the network on which the second disease spreads. If the transmission probability of the second disease is

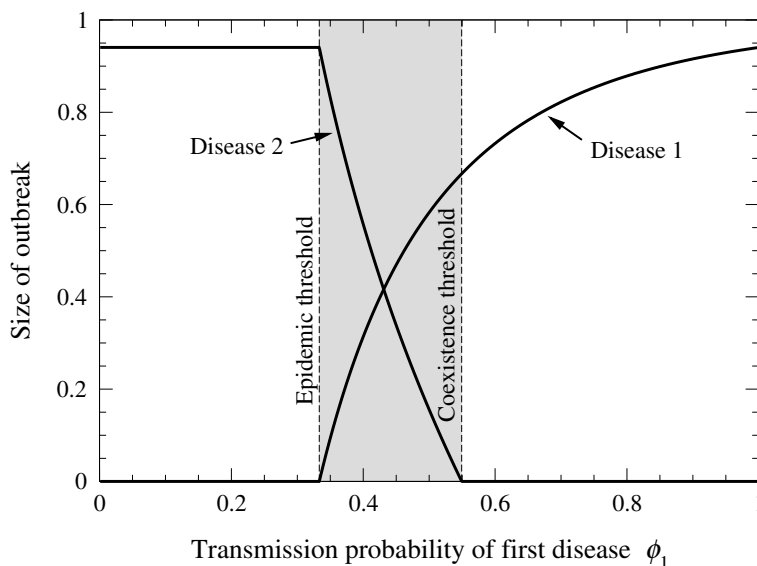


Figure 16.7: Size of outbreaks of two diseases with cross-immunity. The two curves show the sizes of outbreak of two competing diseases as a fraction of total network size on a random graph with Poisson degree distribution and mean degree $c = 3$. The curves are calculated from Eqs. (16.36) and (16.43) while the vertical dashed lines, which represent the positions of the epidemic and coexistence thresholds, are calculated from Eqs. (15.11) and (16.46). After Newman [358].

high enough, then it can spread on this subnetwork and we have “coinfection,” the spread of one disease assisted by another.

We can solve for the behavior of this process on, for example, a configuration model network, but what makes the calculations difficult is that the subnetwork on which the second disease spreads is not, in this case, itself a configuration model. This must be true since the subnetwork in question is, by definition, a percolation cluster of the starting network, and therefore is connected—it has only one component. Since configuration model networks in general have many components, it is clear that this new subnetwork must be something different.

The calculation can nonetheless be done, but it is lengthy and we will not go through all of its intricacies here. The interested reader can find them in Ref. [363]. Figure 16.8 shows the results, again for a Poisson degree distribution with mean degree $c = 3$. The figure shows the sizes of the outbreaks of the two diseases as a function of the transmission probability ϕ_1 of the first disease,

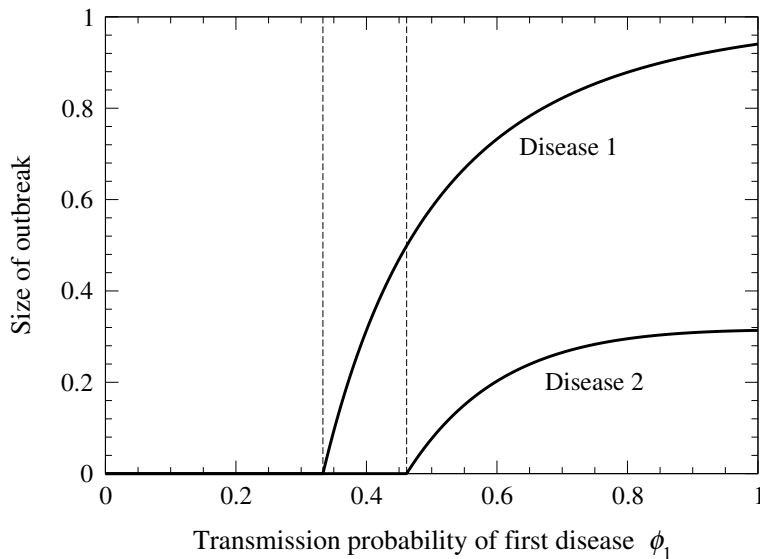


Figure 16.8: Size of outbreaks for coinfection with two diseases. In this calculation, disease 2 can only spread among nodes previously infected with disease 1. The two curves show the sizes of the resulting outbreaks as a fraction of total network size on a random graph with a Poisson degree distribution and mean degree $c = 3$. The transmission probability for disease 2 is held fixed at $\phi_2 = 0.4$ while ϕ_1 is varied from zero to one. The vertical dashed lines represent the positions of the thresholds for each disease to spread. After Newman and Ferrario [363].

with the probability ϕ_2 for the second held fixed. For low values of ϕ_1 (at the left-hand side of the figure) we are below the epidemic threshold for disease 1, so disease 1 does not spread and hence neither does disease 2 (since there are no nodes for it to spread on). As we increase ϕ_1 we eventually pass through the epidemic threshold and disease 1 begins to spread. Once enough individuals are infected with disease 1, disease 2 starts to spread as well, so there are effectively two epidemic thresholds, denoted by the vertical dashed lines in the figure.

16.3.5 COMPLEX CONTAGION

As discussed in Section 16.1.9, the spread of information over networks can be modeled as a contagion-like process, and leads us to consider so-called complex contagions, in which a node only becomes infectious after being infected by two or more others: being infected by only one other is not enough.

Complex contagion is also called *bootstrap percolation* in the physics literature—see Refs. [7, 50, 99, 210].

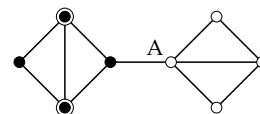
As shown by Baxter *et al.* [50] complex contagions can be treated using percolation techniques similar to those of previous sections. Consider an SIR-style model of complex contagion in which infectious individuals spread the infection to their neighbors with transmission probability ϕ , but a node only becomes infectious itself after it has received the infection from q others.¹² A crucial point to notice about this process is that we can no longer start the outbreak at just a single node. Since the infection only spreads once a node receives it from two or more others, it will never spread at all if it starts at only one node. Indeed, in a large network it is not even adequate to start the disease with a fixed *number* of initial carriers, since the chances of two or more of them being adjacent to the same node vanishes as the size of the network becomes large. In order to get a significant outbreak of a complex contagion, therefore, we must start the infection on a non-vanishing fraction ρ of all nodes. Let us suppose these initial carriers are selected uniformly at random, so that every node has the same probability ρ of having the infection at the start of the outbreak.

Given the identities of the initial carriers we can calculate which other nodes will be infected by a simple iterative procedure. First, the initial carriers spread the infection to their neighbors with probability ϕ and those neighbors become infected if they receive the infection at least q times. Then those newly infected nodes spread the infection to *their* neighbors, who will also become infected if they receive the infection at least q times, including transmission from both the new nodes and the original carriers. And so the process repeats, through as many rounds as are necessary until no more nodes are infected.

It is important to understand that, although this process will correctly tell

¹²Complex contagion has some features in common with the k -cores that we studied in Section 7.2.2. Recall that a k -core in a network is a set of nodes such that each is connected to at least k of the others. The same is also true of the infected individuals in our complex contagion scenario (with k replaced by q), but the two processes are not the same. In part this is because our contagion only spreads with probability ϕ , so there is always a chance that a given node will not get infected. But even if we set $\phi = 1$ there can be differences. The infected individuals in a complex contagion with $\phi = 1$ may constitute a q -core, but more generally they are only a subset of a q -core because of the progressive way the infection spreads, starting from the initial carriers.

Consider the network shown on the right, and suppose we have $q = 2$ and the infection starts from the two circled nodes, spreading with transmission probability $\phi = 1$. The infection will spread to the nodes denoted by the solid circles but will not spread to the four nodes denoted by the open circles on the right, because the node labeled A does not have two infected neighbors. Yet the entire network is a single 2-core, since every node is adjacent to at least two of the others. So in this case the infected nodes do not form a complete 2-core, only a subset.



us who gets infected and who does not, it does not tell us the exact time evolution of the infection or the order in which nodes would get infected in a real outbreak. The rounds of infection described above do not, in general, correspond to real time. It is possible, for instance, for a node in round 3 to actually be infected at an earlier time than a node in round 2. Nonetheless, the iterative process does correctly calculate the total outbreak size and we can use it to develop a percolation-style analytic theory akin to those developed in previous sections for other infection processes. Let us see how the calculation goes using again the example of the configuration model and starting with the simplest case of $q = 2$.

We will denote by u_r the average probability that a node does not receive the infection along a particular one of its edges, on or before round r of the infection process. For general r there are two ways in which the infection can fail to be transmitted along a particular edge. First, the edge might not be occupied, which happens with probability $1 - \phi$. Second, the edge could be occupied (probability ϕ), but the node at its other end is not infected, which can in turn happen in one of two ways: either (1) that node was not an initial carrier and has never received the infection along any of its other edges (that is, along any of its edges excluding the edge we followed to reach it), or (2) it was not an initial carrier and has received the infection along only one edge (which is not enough to become infected). The former happens with probability $(1 - \rho)u_{r-1}^k$, where k is the excess degree of the node, while the latter happens with probability $(1 - \rho)ku_{r-1}^{k-1}(1 - u_{r-1})$. These expressions apply for all values of r including $r = 1$ if we adopt the convention that $u_0 = 1$.

Putting these results together we get a total probability of $1 - \phi + \phi(1 - \rho)u_r^k + \phi(1 - \rho)ku_{r-1}^{k-1}(1 - u_{r-1})$ for the disease to not be transmitted along the edge in question. Averaging over the distribution q_k of the excess degree k , we then arrive at the following equation for u_r :

$$\begin{aligned} u_r &= 1 - \phi + \phi(1 - \rho) \sum_{k=0}^{\infty} q_k u_{r-1}^k + \phi(1 - \rho)(1 - u_{r-1}) \sum_{k=0}^{\infty} q_k k u_{r-1}^{k-1} \\ &= 1 - \phi + \phi(1 - \rho)g_1(u_{r-1}) + \phi(1 - \rho)(1 - u_{r-1})g'_1(u_{r-1}), \end{aligned} \quad (16.47)$$

where g_1 is the generating function for the excess degree distribution, Eq. (12.97), and g'_1 denotes its first derivative, as previously.

In the limit of a large number of rounds the infection will eventually reach every node it is going to infect and u_r will stop changing, so that $u_r = u_{r-1}$ as $r \rightarrow \infty$. For simplicity of notation, let us denote the limiting value of u_r by just u , which represents the average probability that the infection is never transmitted along an edge at any time. Substituting $u_r = u$ and $u_{r-1} = u$ into

Eq. (16.47) then gives

$$u = 1 - \phi + \phi(1 - \rho)g_1(u) + \phi(1 - \rho)(1 - u)g'_1(u). \quad (16.48)$$

If we can solve this equation for u , then the probability that any specific node with degree k is not infected in the limit of long time is equal to the probability that it is not an initial carrier and receives the infection along either zero or one of its k edges but not more, which is $(1 - \rho)u^k + (1 - \rho)(1 - u)u^{k-1}$. Then the probability that a node *is* infected is one minus this expression and, averaging over the degree distribution p_k , the network-wide average probability S of being infected, which is also the fraction of the network infected in the limit of long time, is

$$\begin{aligned} S &= 1 - (1 - \rho) \sum_{k=0}^{\infty} p_k u^k - (1 - \rho)(1 - u) \sum_{k=0}^{\infty} p_k k u^{k-1} \\ &= 1 - (1 - \rho)g_0(u) - (1 - \rho)(1 - u)g'_0(u). \end{aligned} \quad (16.49)$$

Following the same lines of argument it is straightforward to show that for a complex contagion with general q we have

$$u = 1 - \phi + \phi(1 - \rho) \sum_{m=0}^{q-1} \frac{(1 - u)^m}{m!} \frac{d^m g_1}{du^m} \quad (16.50)$$

and

$$S = 1 - (1 - \rho) \sum_{m=0}^{q-1} \frac{(1 - u)^m}{m!} \frac{d^m g_0}{du^m}. \quad (16.51)$$

Between them, Eqs. (16.50) and (16.51) give us our solution for the size S of the outbreak. Note that if we set $q = 1$ and take the limit $\rho \rightarrow 0$, we correctly recover the equations for an ordinary (non-complex) contagion, Eqs. (16.29) and (16.30).

As an example, let us consider again a network with a Poisson degree distribution with mean c , so that $g_0(z) = g_1(z) = e^{c(z-1)}$ and $g'_1(z) = ce^{c(z-1)}$, and examine the case for $q = 2$. Then Eq. (16.48) tells us that

$$u = 1 - \phi + \phi(1 - \rho)[1 + c(1 - u)]e^{c(u-1)}. \quad (16.52)$$

Like Eq. (16.34) for ordinary contagion, this equation has no simple closed-form solution, but we can get a good feel for its behavior using a graphical method. Figure 16.9 shows plots of the right-hand side of the equation as a function of u for several different choices of the parameters, and the point or points at which

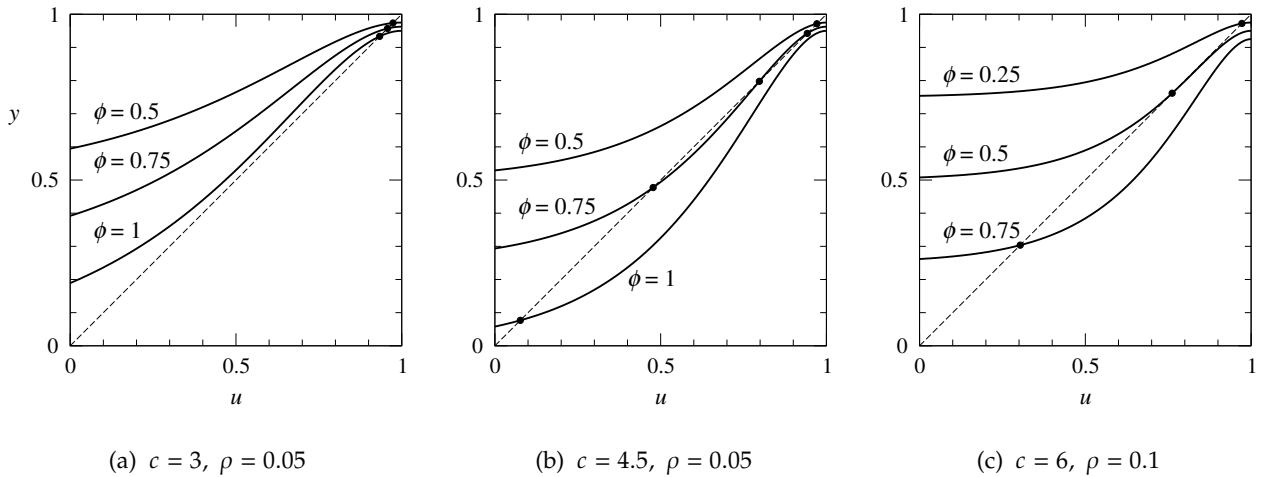


Figure 16.9: Graphical solutions of Eq. (16.52). The curves in these plots show the value of $y = 1 - \phi + \phi(1 - \rho)[1 + c(1 - u)]e^{c(u-1)}$ as a function of u , i.e., the right-hand side of Eq. (16.52), while the dashed diagonal lines are $y = u$. The values of u at which the two cross, marked by the dots, are solutions of (16.52).

the curves cross the dashed diagonal line in each plot represent the solutions for u .

As the figure shows, there are three different types of behavior the system can display, depending on the parameter values. Panel (a) shows the situation when the mean degree of the network is $c = 3$ and the density of initial carriers of the infection is $\rho = 0.05$. Note that the non-zero density of initial carriers implies that the fraction S of infected nodes is always non-zero—in all cases at least a fraction ρ of nodes are infected (in contrast with the epidemics of previous sections, where S was zero below the epidemic threshold).

The three curves in panel (a) of the figure represent the right-hand side of (16.52) for $\phi = 0.5, 0.75$, and 1 . As we can see, for each curve there is just one solution for u , represented by the dots in the top right corner of the plot. Given these solutions we can calculate the size S of the outbreak from Eq. (16.51), which for the Poisson network takes the form

$$S = 1 - (1 - \rho)[1 + c(1 - u)]e^{c(u-1)}. \quad (16.53)$$

Nothing of great interest happens for these parameter values: there is no epidemic threshold, for example, nor sudden behaviors or transitions of any kind. Figure 16.10 shows the full curve of outbreak size S as a function of ϕ

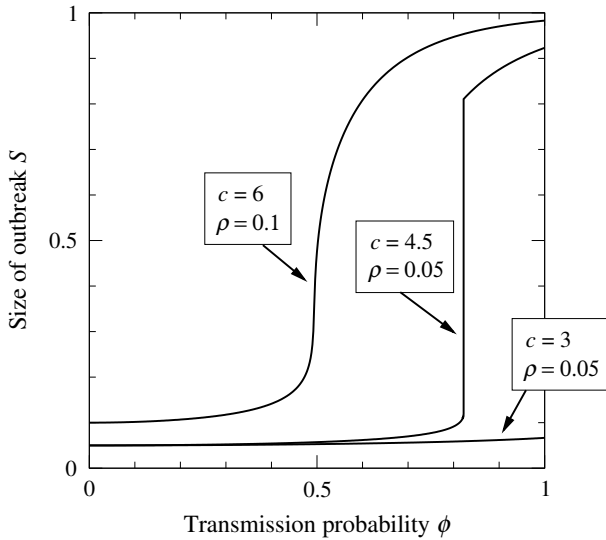


Figure 16.10: Sizes of outbreaks for complex contagions. The three curves show the sizes of outbreaks of a complex contagion with $q = 2$ on a configuration model network with a Poisson degree distribution, for the same choice of parameter values as the three panels of Fig. 16.9. For $c = 3$ and $\rho = 0.05$ (bottom curve) the size of the outbreak is small for all values of ϕ but for $c = 4.5$, $\rho = 0.05$ (middle curve) an avalanche occurs and the outbreak size jumps discontinuously to a much higher value when ϕ is sufficiently large. For $c = 6$, $\rho = 0.1$ (top curve) there is no avalanche, but the outbreak still shows a rapid (though smooth) increase with increasing ϕ .

(the bottom curve in the figure), and the curve is smooth over the whole range of ϕ and the outbreak size remains small—little greater than the baseline value of 0.05 set by the density of initial carriers.

Figure 16.9b shows the situation for parameter values $c = 4.5$ and $\rho = 0.05$, and again the three curves represent $\phi = 0.5, 0.75$, and 1. Now, however, something different happens. For $\phi = 0.5$ there is again just one solution for u , but at $\phi = 0.75$ there are three solutions—the curve crosses and recrosses the diagonal in three different places. Which of these three solutions gives the true behavior of the model? We can answer this question by referring back to the iterative equation (16.47) and recalling that $u_0 = 1$. Starting from this value and iterating we find that u converges to the highest of the three solutions in the figure.

For higher values of ϕ , and in particular for the case of $\phi = 1$ shown in the figure, an interesting thing now happens: the two upper solutions disappear, leaving only the lowest solution of the three. At the point where the upper solutions vanish, the value of u thus *jumps* discontinuously from the highest solution to the lowest. The system shows a *discontinuous phase transition* at this point and the fraction of infected nodes jumps suddenly from a relatively low value to a much higher one. This behavior is clearly visible in Fig. 16.10, which again shows the full curve of S as a function of ϕ (middle curve in the figure), with the discontinuity falling at about $\phi = 0.82$.

This discontinuous phase transition is akin to a kind of epidemic threshold

Discontinuous phase transitions are also sometimes called *first-order* transitions.

for the complex contagion model. Below the phase transition small outbreaks start at many points in the network, everywhere there are two initial carriers close enough together to infect the same third node, but the transmission probability ϕ is too low for these small outbreaks to last long, and they peter out before becoming a true epidemic. The discontinuity represents the point where ϕ becomes large enough that outbreaks no longer peter out. Beyond this point enough new nodes are infected on each round of the complex contagion process to ensure that an equal or larger number get infected on the next round too, and the process becomes self-sustaining. In the limit of large network size it will continue indefinitely. Such runaway infections are sometimes called *avalanches*, by analogy with the runaway dynamics of an avalanche of rocks or snow, where a small initial disturbance produces a self-sustaining landslide.

Figure 16.9c shows a third case, for $c = 6$ and a higher density of initial carriers $\rho = 0.1$, with the three curves in the figure representing solutions for $\phi = 0.25, 0.5$, and 0.75 . For these parameter values there is no discontinuity in S . There is no value of ϕ at which the curve crosses the diagonal more than once, so the value of S is a smooth function of ϕ , as shown in the top curve in Fig 16.10. However, the value does change quite rapidly from low to high around $\phi = 0.5$. It's as if the system is on the verge of undergoing an avalanche, but in the end it doesn't quite get there.

16.4 TIME-DEPENDENT PROPERTIES OF EPIDEMICS ON NETWORKS

The techniques of the previous sections can tell us about the late-time properties of infections on networks, such as how many people will eventually be affected in an outbreak of a disease. They do not, however, tell us about the detailed progression of an outbreak over time. If we want to know about the temporal evolution of infections then we need another approach. Moreover, the techniques we have used so far cannot tell us about even the late-time behavior of models with reinfection, such as the SIS and SIRS models of Sections 16.1.5 and 16.1.6. For these models the equivalence between epidemics and percolation that we used in previous sections does not hold, and to understand their behavior, including at long times, we need to address the dynamics of the epidemic.

There are a number of approaches for calculating the dynamics of epidemics on networks, some exact and some approximate. Given a specific network one can always perform computer simulations and get numerical answers for typical disease outbreaks. Analytic approaches, on the other hand, offer more insight, but they are mostly confined to specific model networks, such as random graphs and their generalizations. In the following sections we will look at some

of the most straightforward and general approaches to epidemic dynamics on networks, starting with the simple SI model and progressing to more complex (and interesting) models in later sections.

16.5 TIME-DEPENDENT PROPERTIES OF THE SI MODEL

The analytic treatment of the time-dependent properties of epidemic models centers on the probabilities for nodes to be in specific disease states at specific times. One can imagine having repeated outbreaks of the same disease on the same network, starting from the same initial conditions, and calculating, for example, the average probabilities $s_i(t)$ and $x_i(t)$ that node i is susceptible or infective at time t . Given the adjacency matrix of a network one can write down differential equations for the evolution of such quantities similar to those for the fully mixed models of Section 16.1. Consider for instance the SI model.

As we have seen, an SI outbreak, starting at a single randomly chosen node, eventually spreads to all members of the component containing that node. Our main interest is in epidemics occurring in the giant component of the network, since all other outbreaks will only affect a small component and then die out, so let us focus on the giant component.

Consider a given node i . If the node is not a member of the giant component then by hypothesis $s_i = 1$ at all times—the node is always susceptible since we are assuming the epidemic to take place in the giant component. If i is in the giant component then we can write down a differential equation for s_i by considering the probability that i becomes infected between times t and $t + dt$. To become infected, an individual must catch the disease from a neighboring individual j , meaning j must already be infected, which happens with probability $x_j = 1 - s_j$, and must transmit the disease during the given time interval, which happens with probability βdt . In addition we also require that i be susceptible in the first place, which happens with probability s_i . Multiplying these probabilities and then summing over all neighbors of i , the total probability of i becoming infected is $\beta s_i dt \sum_j A_{ij} x_j$, where A_{ij} is an element of the adjacency matrix. Thus s_i obeys the differential equation

$$\frac{ds_i}{dt} = -\beta s_i \sum_j A_{ij} x_j = -\beta s_i \sum_j A_{ij} (1 - s_j). \quad (16.54)$$

Note the leading minus sign on the right-hand side—the probability of being susceptible goes down when nodes become infected.

Similarly we can write an equation for x_i thus:

$$\frac{dx_i}{dt} = \beta s_i \sum_j A_{ij} x_j = \beta(1 - x_i) \sum_j A_{ij} x_j, \quad (16.55)$$

although the two equations, (16.54) and (16.55), are really the same equation, related to one another by $s_i + x_i = 1$.

We will use the same initial conditions as we did in the fully mixed case, assuming that the disease starts with either a single infected node or a small number c of nodes, chosen uniformly at random, so that $x_i = c/n$ and $s_i = 1 - c/n$ for all i . In the limit of large system size n , these become $x_i = 0$, $s_i = 1$, and we will use this large- n limit to simplify some of the expressions derived in this and the following sections.¹³

Equation (16.54) is not solvable in closed form for general A_{ij} but we can calculate some features of its behavior by considering suitable limits. Consider, for example, the behavior of the system at early times. For large n , and assuming initial conditions as above, x_i will be small in this regime. Working with Eq. (16.55) and ignoring terms of quadratic order in small quantities, we have

$$\frac{dx_i}{dt} = \beta \sum_j A_{ij} x_j, \quad (16.56)$$

or in matrix form

$$\frac{d\mathbf{x}}{dt} = \beta \mathbf{A}\mathbf{x}, \quad (16.57)$$

where \mathbf{x} is the vector with elements x_i .

Now let us write \mathbf{x} as a linear combination of the eigenvectors of the adjacency matrix:

$$\mathbf{x}(t) = \sum_{r=1}^n a_r(t) \mathbf{v}_r, \quad (16.58)$$

where \mathbf{v}_r is the eigenvector with eigenvalue κ_r and $a_r(t)$ is a time-dependent coefficient. Then

$$\sum_{r=1}^n \frac{da_r}{dt} \mathbf{v}_r = \frac{d\mathbf{x}}{dt} = \beta \mathbf{A}\mathbf{x} = \beta \mathbf{A} \sum_{r=1}^n a_r(t) \mathbf{v}_r = \beta \sum_{r=1}^n \kappa_r a_r(t) \mathbf{v}_r. \quad (16.59)$$

¹³If we choose the initial carriers of the disease uniformly from the entire network in this way, then there is a chance they will fall outside the giant component. If, as discussed earlier, we want to be certain that an outbreak takes place in the giant component, then we can choose the initial carriers randomly among only the nodes in that giant component. In this case, the total number of nodes n in the expressions above would be replaced by the number in the giant component.

Thus, comparing terms in \mathbf{v}_r , we have

$$\frac{da_r}{dt} = \beta \kappa_r a_r, \quad (16.60)$$

which has the solution

$$a_r(t) = a_r(0) e^{\beta \kappa_r t}. \quad (16.61)$$

Substituting this expression back into Eq. (16.58) we get our solution for the early-time behavior of the epidemic:

$$\mathbf{x}(t) = \sum_{r=1}^n a_r(0) e^{\beta \kappa_r t} \mathbf{v}_r. \quad (16.62)$$

The fastest growing term in this expression is the term corresponding to the largest eigenvalue κ_1 . Assuming this term dominates over the others we will get

$$\mathbf{x}(t) \sim e^{\beta \kappa_1 t} \mathbf{v}_1. \quad (16.63)$$

So we expect the number of infected individuals to grow exponentially at short times, just as it does in the fully mixed version of the SI model (see Section 16.1.1 and Fig. 16.1), but now with an exponential constant that depends not just on β but also on the leading eigenvalue of the adjacency matrix. Moreover, the probability of infection in this early period varies from node to node as the corresponding element of the leading eigenvector \mathbf{v}_1 . The elements of the leading eigenvector of the adjacency matrix are the same quantities that in other circumstances we called the eigenvector centrality—see Section 7.1.2. Thus eigenvector centrality is an approximate measure of the probability of early infection of a node in an SI epidemic.

At long times in the SI model the probability of infection of any node in the giant component tends to one (again assuming that the epidemic takes place in the giant component). Thus, overall we expect the SI epidemic to have a similar form to that seen in the fully mixed version of the model, producing curves qualitatively like that in Fig. 16.1 but with nodes of higher eigenvector centrality becoming infected faster than those of lower.

Although Eqs. (16.54) and (16.55) and their derivation seem reasonable they are not, in fact, precisely correct, as we can see by solving the equations numerically. Figure 16.11a shows the results of such a numerical solution (the curve labeled “first-order”) on a configuration model network, compared against an average of a large number of epidemics with the same β simulated directly on the same network (the circular dots). As the figure shows, the agreement between the two is good, but definitely not perfect.

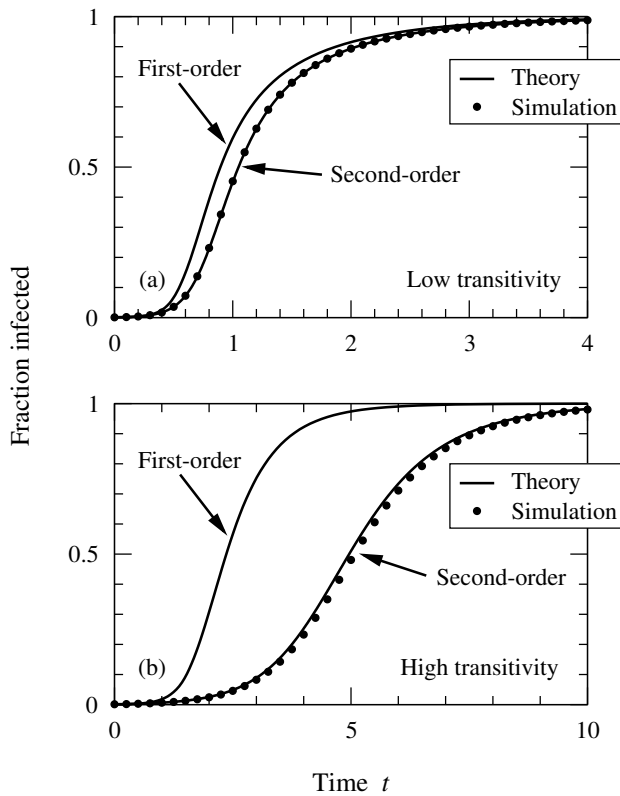


Figure 16.11: Comparison of theory and simulation for the SI model on two different networks. (a) The fraction of infected individuals as a function of time on the giant component of a network with low transitivity (i.e., low clustering coefficient), calculated by numerical solution of the differential equations for the first- and second-order moment closure methods, and by direct simulation. (b) The same comparison for a network with high transitivity. The networks have one million nodes each and the transmission rate is $\beta = 1$ in all cases. Simulation results are averaged over 500 runs.

The reason for this disagreement is an interesting one. Equation (16.54) may appear to be a straightforward generalization of the equivalent equation for the fully mixed SI model, Eq. (16.4), but there are some subtleties involved. The right-hand side of the equation contains two average quantities, s_i and x_j , and in multiplying these quantities we are implicitly assuming that the product of the averages is equal to the average of their product. In the fully mixed model this is true (for large n) because of the mixing itself, but in the present case it is, in general, not, because the probabilities are not independent. The quantity s_i measures a node's probability of being susceptible and x_j measures the probability of its neighbor being infected. It should come as no surprise that in general these quantities will be correlated between neighboring nodes. Correlations of this type can be incorporated into our calculations, at least approximately, by using a so-called pair approximation or moment closure method, as described in the following section.

16.5.1 PAIR APPROXIMATION

Correlations between the disease states of different nodes can be handled by augmenting our theory to take account of the joint probabilities for pairs of nodes to have given pairs of states. To handle such joint probabilities we will need to make our notation a little more sophisticated. Let us denote by $\langle s_i \rangle$ the average probability that node i is susceptible. This is the same quantity that we previously called s_i , but, as we will see, it will be useful to indicate the average explicitly with angle brackets $\langle \dots \rangle$. You can think of $s_i(t)$ as now being a variable with value one if i is susceptible at time t and zero otherwise and $\langle s_i \rangle$ as being the average of this quantity over many different instances of disease outbreaks on the same network. Similarly $\langle x_i \rangle$ will be the average probability that i is infected. And $\langle s_i x_j \rangle$ indicates the average probability that i is susceptible *and* j is infected at the same time.

In this notation it is straightforward to write down a truly exact version of Eq. (16.54), taking correlations into account. It is

$$\frac{d\langle s_i \rangle}{dt} = -\beta \sum_j A_{ij} \langle s_i x_j \rangle. \quad (16.64)$$

Equation (16.54) is an approximation to this true equation in which we assume that $\langle s_i x_j \rangle \simeq \langle s_i \rangle \langle x_j \rangle$.

The trouble with Eq. (16.64) is that we cannot solve it directly because it contains the unknown quantity $\langle s_i x_j \rangle$ on the right-hand side. To find this quantity, we need another equation for $\langle s_i x_j \rangle$, which we can deduce as follows.

To reach the state in which i is susceptible and j is infected in an SI model it must be the case that both i and j are susceptible to begin with and then j becomes infected. Even though i and j are neighbors j cannot catch the disease from i since i is not infected, so j must catch the disease from some other neighboring node k , which itself must be infected. In our new notation, the probability for the configuration in which i and j are susceptible and k is infected is $\langle s_i s_j x_k \rangle$. If we have this configuration, then the disease will be transmitted from k to j with rate β . Summing over all neighbors k except for i , the total rate at which j becomes infected is then $\beta \sum_{k(\neq i)} A_{jk} \langle s_i s_j x_k \rangle$.

Unfortunately, this is not the end of the story because $\langle s_i x_j \rangle$ can also *decrease*—it decreases if i becomes infected. This can happen in two different ways. Either i can catch the disease from its infected neighbor j , which happens with rate $\beta \langle s_i x_j \rangle$, or it can catch it from another infected neighbor $l \neq j$, which happens with rate $\beta \langle x_l s_i x_j \rangle$. Summing the latter expression over all neighbors l other than j gives a total rate of $\beta \sum_{l(\neq j)} A_{il} \langle x_l s_i x_j \rangle$.

Putting all of these terms together, with minus signs for those that decrease

the probability, we get an equation for $\langle s_i x_j \rangle$ thus:

$$\frac{d\langle s_i x_j \rangle}{dt} = \beta \sum_{k(\neq i)} A_{jk} \langle s_i s_j x_k \rangle - \beta \sum_{l(\neq j)} A_{il} \langle x_l s_i x_j \rangle - \beta \langle s_i x_j \rangle. \quad (16.65)$$

In theory this equation will now allow us to calculate $\langle s_i x_j \rangle$. In practice, however, it involves yet more terms that we don't know on the right-hand side, the three-variable averages $\langle s_i s_j x_k \rangle$ and $\langle x_l s_i x_j \rangle$. We can write down further equations for these averages but, as you might guess, those equations involve still higher-order (four-variable) terms, and so forth. The succession of equations will never end—in the language of mathematics, it doesn't *close*—and so it looks as though it will be of no use to us.¹⁴

In fact, however, we can still make progress by approximating our three-variable averages with appropriate combinations of one- and two-variable averages, which allows us to close the equations and get a set we can actually solve. This type of approach is called a *moment closure method*. The moment closure method at the level of two-variable averages that we discuss here is also called a *pair approximation*.

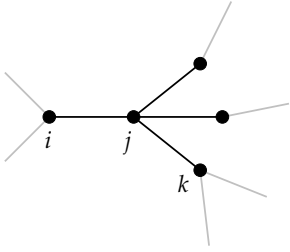
In fact, our first attempt at writing equations for the SI model on a network, Eq. (16.54), was itself a simple moment-closure method. We approximated the true equation, Eq. (16.64), by writing $\langle s_i x_j \rangle \simeq \langle s_i \rangle \langle x_j \rangle$, closing the equations at the level of one-variable averages. Going a step further and closing at the level of two-variable averages makes our equations more precise, and in fact, as we will see, this “second-order” moment closure is exact for some networks, although only approximate for others. Even in the latter case, however, the method gives a remarkably good approximation. The approximation can be made better still by going to third order, but the equations rapidly become more complicated and researchers have rarely used moment closure methods beyond the second-order, pair-approximation level.

The pair approximation is relatively straightforward however. Starting with Eq. (16.65) our goal is to approximate the three-variable averages on the right-hand side with lower-order ones. To do this we first write

$$\langle s_i s_j x_k \rangle = P(i, j \in S, k \in I) = P(i, j \in S)P(k \in I | i, j \in S), \quad (16.66)$$

where $P(i \in S)$ means the probability that node i is in the set S of susceptible nodes. We know that i and j are neighbors in the network and that j and k

¹⁴On a finite network with n nodes the equations will in fact close once we get all the way up to combinations of n variables, but this limit is not useful in practice as the equations will become unmanageably numerous and complicated long before we reach it.



The pair approximation is exact if the only connection between nodes i and k is through node j .

are neighbors. Our approximation involves assuming that the disease state of k does not depend on the disease state of i . This is a good approximation, indeed not an approximation at all, if the only path in the network from i to k is through j . In that case, given that we know j to be susceptible, there is no way that the disease state of i can affect that of k because there is no way the disease could have spread from i to k . On the other hand, if there is another path from i to k that avoids node j , then the disease could take that path, which will introduce correlations between i and k : if i is infected then k is more likely to be infected too. In that case our approximation is just that—an approximation—although as we will see it may be a very good one.

Assuming the state of k to be independent of the state of i , we have

$$P(k \in I | i, j \in S) = P(k \in I | j \in S) = \frac{P(j \in S, k \in I)}{P(j \in S)} = \frac{\langle s_j x_k \rangle}{\langle s_j \rangle}. \quad (16.67)$$

Putting Eqs. (16.66) and (16.67) together, we then have

$$\langle s_i s_j x_k \rangle = \frac{\langle s_i s_j \rangle \langle s_j x_k \rangle}{\langle s_j \rangle}. \quad (16.68)$$

We can write a similar expression for the other three-variable average appearing in Eq. (16.65):

$$\langle x_l s_i x_j \rangle = \frac{\langle x_l s_i \rangle \langle s_i x_j \rangle}{\langle s_i \rangle}, \quad (16.69)$$

and, substituting both into Eq. (16.65), we then get the pair approximation equation

$$\frac{d\langle s_i x_j \rangle}{dt} = \beta \frac{\langle s_i s_j \rangle}{\langle s_j \rangle} \sum_{k(\neq i)} A_{jk} \langle s_j x_k \rangle - \beta \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \sum_{l(\neq j)} A_{il} \langle s_i x_l \rangle - \beta \langle s_i x_j \rangle. \quad (16.70)$$

This equation now contains only averages over two variables at a time. It does also contain a new average $\langle s_i s_j \rangle$ that we have not encountered before, but this can easily be rewritten as $\langle s_i s_j \rangle = \langle s_i (1 - x_j) \rangle = \langle s_i \rangle - \langle s_i x_j \rangle$ and so our equation becomes

$$\frac{d\langle s_i x_j \rangle}{dt} = \beta \frac{\langle s_i \rangle - \langle s_i x_j \rangle}{\langle s_j \rangle} \sum_{k(\neq i)} A_{jk} \langle s_j x_k \rangle - \beta \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \sum_{l(\neq j)} A_{il} \langle s_i x_l \rangle - \beta \langle s_i x_j \rangle. \quad (16.71)$$

This equation is more complex than Eq. (16.54) but it can be simplified by rewriting it as follows. Let us define p_{ij} to be the conditional probability that j is infected given that i is not:

$$p_{ij} = P(j \in I | i \in S) = \frac{P(i \in S, j \in I)}{P(i \in S)} = \frac{\langle s_i x_j \rangle}{\langle s_i \rangle}. \quad (16.72)$$

Then the time evolution of p_{ij} is given by

$$\begin{aligned}
 \frac{dp_{ij}}{dt} &= \frac{d}{dt} \left(\frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \right) \\
 &= \frac{1}{\langle s_i \rangle} \frac{d\langle s_i x_j \rangle}{dt} - \frac{\langle s_i x_j \rangle}{\langle s_i \rangle^2} \frac{d\langle s_i \rangle}{dt} \\
 &= \beta \left(1 - \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \right) \sum_{k(\neq i)} A_{jk} \frac{\langle s_j x_k \rangle}{\langle s_j \rangle} - \beta \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \sum_{l(\neq j)} A_{il} \frac{\langle s_i x_l \rangle}{\langle s_i \rangle} \\
 &\quad - \beta \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} + \beta \frac{\langle s_i x_j \rangle}{\langle s_i \rangle} \sum_l A_{il} \frac{\langle s_i x_l \rangle}{\langle s_i \rangle} \\
 &= \beta(1 - p_{ij}) \sum_{k(\neq i)} A_{jk} p_{jk} - \beta p_{ij} \sum_{l(\neq j)} A_{il} p_{il} - \beta p_{ij} + \beta p_{ij} \sum_l A_{il} p_{il}, \quad (16.73)
 \end{aligned}$$

where we have used Eqs. (16.64) and (16.71) in the third line. All but one of the terms in the two sums over l now cancel out, leaving us with the relatively simple equation

$$\frac{dp_{ij}}{dt} = \beta(1 - p_{ij}) \left[-p_{ij} + \sum_{k(\neq i)} A_{jk} p_{jk} \right], \quad (16.74)$$

where we have used the fact that $A_{ij} = 1$ (since i and j are neighbors). We can also rewrite Eq. (16.64) in terms of p_{ij} thus:

$$\frac{d\langle s_i \rangle}{dt} = -\beta \langle s_i \rangle \sum_j A_{ij} p_{ij}, \quad (16.75)$$

which has the solution

$$\langle s_i(t) \rangle = \langle s_i(0) \rangle \exp \left(-\beta \sum_j A_{ij} \int_0^t p_{ij}(t') dt' \right). \quad (16.76)$$

If we can solve Eq. (16.74) for p_{ij} then we can substitute the result into Eq. (16.76) to get our solution for the time evolution of the epidemic. Note that there are two equations of the form (16.74) for each edge in the network, since p_{ij} is not symmetric in i and j .

Figure 16.11a on page 649 shows results from a numerical solution of these equations (the curve marked “second-order”), again on a configuration model network and, as the figure shows, the calculation now agrees very well with the simulation results represented by the dots in the figure. By accounting for correlations between adjacent nodes we have created a much more accurate theory.

This near-perfect agreement, however, is something of a special case. Configuration model networks are locally tree-like, meaning they have no short loops (see Section 12.4), and, as discussed earlier, our second-order moment closure approximation is exact when non-adjacent nodes i and k have only a single path between them through some intermediate node j . When there are no short loops in our network this is true to an excellent approximation—the only other way to get from i to k in such a network is by going around a long loop and the length of such loops dilutes any resulting correlations between the states of i and k , often to the point where they can be ignored. The network used in the simulations for Fig. 16.11a was sufficiently large (a million nodes) and the resulting loops were sufficiently long that the pair approximation equations are an excellent approximation in this case, which is why the agreement is so good in the figure.

Unfortunately, as we saw in Section 7.3, most real social networks have a lot of short loops, which raises the question of how well our method does on such networks. Figure 16.11b shows a comparison between the predictions of our equations and direct simulations for a network with many short loops,¹⁵ for both the simple first-order moment closure method of Eq. (16.54) and for our more sophisticated second-order approach. As the plot shows, the first-order calculation agrees quite poorly with the simulations, its predictions being inaccurate enough to be of little use in this case. The second-order equations, however, still do remarkably well. Their predictions are not in perfect agreement with the simulations, but they are close.

Thus, the pair approximation method offers a significant improvement for networks both with and without short loops, providing a usefully accurate approximation in the former case and being essentially exact in the latter.

16.5.2 DEGREE-BASED APPROXIMATION FOR THE SI MODEL

The analysis of the previous section gives asymptotically exact equations for the dynamics of the SI model on a network with no short loops and an excellent approximation in other cases. Unfortunately, the equations cannot, in general, be solved analytically, even for simple networks such as those generated by the configuration model. (The results in Fig. 16.11 were derived by integrating the equations numerically.)

In this section we describe an alternative approximate approach that gives good, though not perfect, results in practice and produces equations that can be solved analytically. Moreover, the method can, as we will see, be generalized

¹⁵The network was generated using the clustered network model of Ref. [353].

to other epidemic models such as the SIR model. The method was pioneered by Pastor-Satorras and co-workers [47, 48, 382, 383], though it has precursors in earlier work by May and others [306, 325]. It takes its simplest form when applied to networks drawn from the configuration model and so it is on this model that we focus here, although in principle the method can be extended to other networks.

Consider a disease propagating on a configuration model network, i.e., a random graph with a given degree distribution p_k , as discussed in Chapter 12. As before we focus on outbreaks taking place in the giant component of the network, this being the case of most interest—outbreaks in small components by definition die out quickly and do not give rise to epidemics.

An important point to note is that the degree distribution of nodes in the giant component of a configuration model network is not the same as the degree distribution of nodes in the network as a whole. The probability of a node belonging to the giant component goes up with degree because more edges means more chances to be connected to the giant component (see Section 12.6). This means that the degree distribution in the giant component is skewed towards higher degrees. (For a start, note that there are trivially no nodes of degree zero in the giant component, since by definition such nodes are not attached to any others.) In this section we will, as previously, denote the degree distribution and the excess degree distribution in our calculations by p_k and q_k , but bear in mind that these are for nodes in the giant component, which means they are not the same as the distributions for the network as a whole.

The approximation introduced by Pastor-Satorras *et al.* was to assume that all nodes of the same degree have the same probability of infection at any given time. Certainly this *is* an approximation. The probability of infection for a node of degree, say, five situated in the middle of the dense core of a network will very likely be larger than the probability for a node of degree five that is out on the periphery. Nonetheless, if the distribution of probabilities for nodes of given degree is relatively narrow, it may be a good approximation to set them all equal to the same value. And in practice, as we have said, the approximation appears to work well.

Returning, for the sake of simplicity, to our earlier notation style, let us define $s_k(t)$ and $x_k(t)$ to be the probabilities that a node with degree k is susceptible or infected, respectively, at time t . Now consider a susceptible node i . To become infected, i has to contract the infection from one of its network neighbors. The probability that a particular neighbor j is infected depends on the neighbor's degree, but we must be careful. By hypothesis node i is not infected and so j cannot have caught the disease from i . If j is infected it must have caught the disease from one of its other neighbors. In effect this reduces the degree of j by

one— j will have the same probability of being infected at the current time as the average node with degree one less. To put that another way, j 's probability of infection depends upon its excess degree, the number of edges it has other than the edge we followed from i to reach it. Node j 's probability of infection is thus x_k , but where k indicates the excess degree, not the total degree.

The advantage of the degree-based approach now becomes clear: the probability of j being infected depends, in this approach, only on j 's excess degree and not on i 's degree. By contrast, the conditional probability p_{ij} in the pair approximation of Section 16.5.1 was a function of two indices, making the equations more complicated.

To derive the equations for the degree-based approximation, consider the probability that node i becomes infected between times t and $t + dt$. To become infected it must catch the disease from one of its neighbors, meaning that the neighbor must be infected. The probability of a neighbor being infected is x_k where k is the excess degree of the neighbor, and the excess degree is distributed according to the distribution q_k of Eq. (12.16), which means that the average probability that the neighbor is infected is

$$v(t) = \sum_{k=0}^{\infty} q_k x_k(t). \quad (16.77)$$

If the neighbor is infected, then the probability that the disease will be transmitted to node i in the given time interval is βdt . Then the overall probability of transmission from a single neighbor during the time interval is $\beta v(t) dt$ and the probability of transmission from any neighbor is $\beta k v(t) dt$, where k is now the number of i 's neighbors, its total degree. In addition we also require that i itself be susceptible, which happens with probability $s_k(t)$, so our final probability that i becomes infected is $\beta k v s_k dt$. Thus the rate of change of s_k is given by

$$\frac{ds_k}{dt} = -\beta k v s_k. \quad (16.78)$$

This equation can be solved exactly. We can formally integrate it thus:

$$s_k(t) = s_0 \exp\left(-\beta k \int_0^t v(t') dt'\right), \quad (16.79)$$

where we have fixed the integration constant so that all nodes have probability s_0 of being susceptible at $t = 0$.

Although we don't yet know the form of the function $v(t)$, Eq. (16.79) tells us that s_k depends on k as a simple power of some universal k -independent function $u(t)$:

$$s_k(t) = s_0 [u(t)]^k, \quad (16.80)$$

where

$$u(t) = \exp\left(-\beta \int_0^t v(t') dt'\right). \quad (16.81)$$

Rather than evaluating the integral explicitly, it's easier to calculate $u(t)$ by differentiating to get

$$\frac{du}{dt} = -\beta uv, \quad (16.82)$$

and evaluating $v(t)$ by making the substitution $x_k = 1 - s_k$ in Eq. (16.77) to get

$$v(t) = \sum_{k=0}^{\infty} q_k(1 - s_k) = \sum_{k=0}^{\infty} q_k(1 - s_0 u^k) = 1 - s_0 g_1(u), \quad (16.83)$$

where $g_1(u)$ is the generating function for q_k and we have made use of Eq. (16.80) and the fact that $\sum_k q_k = 1$. Substituting this expression into Eq. (16.82) then gives us

$$\frac{du}{dt} = -\beta u [1 - s_0 g_1(u)], \quad (16.84)$$

which is a straightforward linear differential equation for u that, given the degree distribution, can be solved by direct integration with the initial condition $u = 1$ at $t = 0$ (which comes from setting $t = 0$ in Eq. (16.81)).

Once we have $u(t)$, we can calculate s_k from Eq. (16.80). Then we can calculate the total fraction $x(t)$ of infected individuals in the network by averaging over k thus:

$$x(t) = \sum_{k=1}^{\infty} p_k x_k(t) = \sum_{k=1}^{\infty} p_k (1 - s_0 u^k) = 1 - s_0 g_0(u). \quad (16.85)$$

Note that the sums here start at $k = 1$ because there are no nodes of degree zero in the giant component.

Equations (16.84) and (16.85) between them give us an approximate solution for the SI model on the giant component of a configuration model network with any degree distribution.

Although this approach is elegant in principle, in most practical cases we cannot solve Eq. (16.84) in closed form. Even without finding a solution, however, we can see the basic form it must take. First of all, as we have said, $u = 1$ at $t = 0$ by Eq. (16.81). And since $v(t)$ represents the average probability that your neighbor is infected, it is, by definition, positive and non-decreasing with time, in which case Eq. (16.81) also implies that $u(t)$ always decreases and tends to zero as $t \rightarrow \infty$. This in turn implies that at long times Eq. (16.84) becomes

$$\frac{du}{dt} = -\beta u [1 - s_0 g_1(0)] = -\beta u (1 - s_0 p_1 / \langle k \rangle), \quad (16.86)$$

where we have neglected terms of order u^2 and smaller. Hence $u(t)$ decays exponentially as $e^{-\beta(1-s_0 p_1/\langle k \rangle)t}$ when $t \rightarrow \infty$. Assuming the infection starts with only one or a handful of cases, so that $s_0 = 1 - c/n$ for some constant c , we have $s_0 \rightarrow 1$ in the limit of large n and

$$u(t) \sim e^{-\beta(1-p_1/\langle k \rangle)t}. \quad (16.87)$$

Note that the long-time behavior depends on the fraction p_1 of nodes with total degree one. This is because these are the last nodes to be infected—individuals with only one contact are best protected from infection, although even they are guaranteed to become infected in the end. In networks where the fraction p_1 is zero or very small we have $u(t) \sim e^{-\beta t}$ and the functional form of the long-time behavior depends only on the infection rate and not on the network structure.

At short times we can write $u = 1 - \epsilon$ and to leading order Eq. (16.84) becomes

$$\frac{d\epsilon}{dt} = \beta [x_0 + g'_1(1)\epsilon], \quad (16.88)$$

where $x_0 = 1 - s_0$ is the initial value of x .¹⁶ This has the solution

$$\epsilon(t) = \frac{x_0}{g'_1(1)} [e^{\beta g'_1(1)t} - 1], \quad (16.89)$$

where we have made use of the initial condition $\epsilon = 0$. Equivalently, we can write

$$u(t) = 1 - \epsilon = 1 - \frac{x_0}{g'_1(1)} [e^{\beta g'_1(1)t} - 1]. \quad (16.90)$$

Given the short- and long-time behavior and the fact that $u(t)$ is monotonically decreasing, we can now guess that $u(t)$ has a form something like that of Fig. 16.12. Then, since g_0 is a monotonically increasing function of its argument, $x(t)$ in Eq. (16.85) has a similar shape but turned upside down, so that it looks qualitatively

similar to the curve for the fully mixed version of the model shown in Fig. 16.1, although quantitatively it may be different.

The initial growth of $x(t)$ can be calculated by putting $u = 1 - \epsilon$ in Eq. (16.85) to give $g_0(1 - \epsilon) \simeq 1 - g'_0(1)\epsilon$ and

$$x(t) = 1 - s_0 + s_0 g'_0(1)\epsilon = x_0 \left[1 + \frac{g'_0(1)}{g'_1(1)} (e^{\beta g'_1(1)t} - 1) \right], \quad (16.91)$$

¹⁶As previously, we have assumed that x_0 is small, and we keep terms to leading order only in both ϵ and x_0 ; i.e., we have dropped terms going as $x_0\epsilon$.

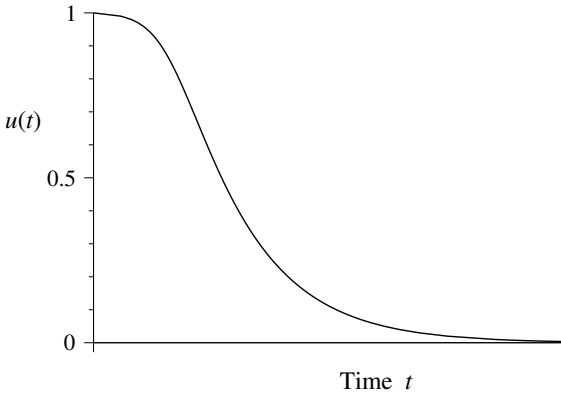


Figure 16.12: The function $u(t)$ in the solution of the SI model. Generically we expect $u(t)$ to have the form sketched here: it is monotonically decreasing from an initial value of 1 and has an exponential tail at long times.

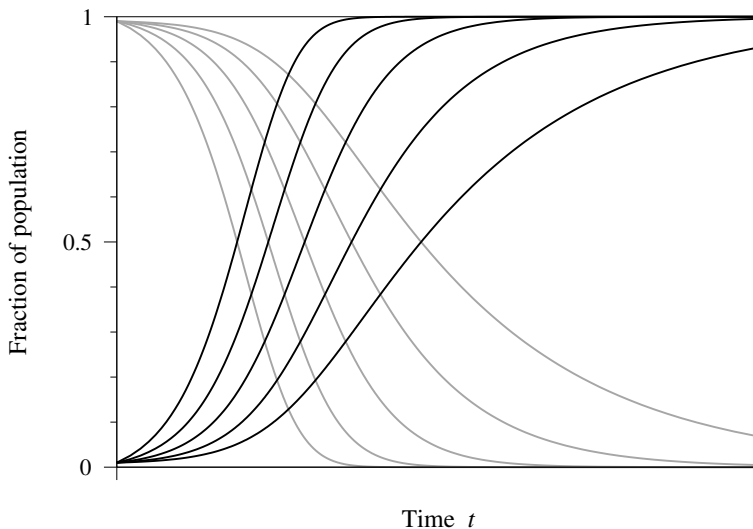


Figure 16.13: Fractions of susceptible and infected nodes of various degrees in the SI model. The various curves show the fraction of nodes of degree k that are susceptible (gray) and infected (black) as a function of time for $k = 1, 2, 4, 8,$ and 16 . The highest values of k give the fastest changing (leftmost) curves and the lowest values the slowest changing. The curves were calculated by integrating Eq. (16.84) numerically with $\beta = 1$ and a Poisson degree distribution with mean degree four.

where we have again set $s_0 = 1$. Thus, as we would expect, the initial growth of infection is roughly exponential.

Note the appearance of $g'_1(1)$ in the exponential in Eq. (16.91). As we saw in Eq. (12.110) on page 410, $g'_1(1)$ is equal to the ratio c_2/c_1 of the average number of second neighbors of a node to the average number of first neighbors, and hence is a measure of how fast the network branches as we move away from a node. It should be not surprising, therefore, to see that this same quantity—along with the transmission rate β —controls the rate at which the disease spreads in our SI model.

Another interesting feature of the model is the behavior of the quantities $s_k(t)$ that measure the probability that a node of a given degree is susceptible. Since these quantities are all proportional to powers of $u(t)$ —see Eq. (16.80)—they form a family of curves as shown in Fig. 16.13. Thus, as we might expect, the nodes with highest degree are the ones that become infected first, on average, while those with low degree hold out longer.

16.6 TIME-DEPENDENT PROPERTIES OF THE SIR MODEL

It is relatively straightforward to extend the techniques of the previous sections to the more complex (and interesting) SIR model. Again we concentrate on outbreaks taking place in the giant component of the network and we define $s_i(t)$, $x_i(t)$, and $r_i(t)$ to be the probabilities that node i is susceptible, infected, or recovered respectively at time t . The equivalent of Eqs. (16.54) and (16.55)—the first-order moment-closure approximation in which all probabilities are assumed independent—are straightforward to write down. The evolution of s_i is (approximately) governed by the same equation as before:

$$\frac{ds_i}{dt} = -\beta s_i \sum_j A_{ij} x_j, \quad (16.92)$$

while x_i and r_i obey

$$\frac{dx_i}{dt} = \beta s_i \sum_j A_{ij} x_j - \gamma x_i, \quad (16.93)$$

$$\frac{dr_i}{dt} = \gamma x_i, \quad (16.94)$$

where, as previously, γ is the recovery rate, i.e., the probability per unit time that an infected individual will recover.¹⁷

We can choose the initial conditions in various ways, but let us here make the same assumption as we did for the SI model, that at $t = 0$ we have a small number c of infected individuals and everyone else is susceptible, so that $s_i(0) = 1 - c/n$, $x_i(0) = c/n$, and $r_i(0) = 0$.

We cannot solve Eqs. (16.92) to (16.94) exactly, but we can extract some useful results by examining their behavior at early times. In the limit $t \rightarrow 0$, x_i is small and $s_i = 1 - c/n$, which tends to 1 as n becomes large, so Eq. (16.93) can be approximated as

$$\frac{dx_i}{dt} = \beta \sum_j A_{ij} x_j - \gamma x_i = \sum_j (\beta A_{ij} - \gamma \delta_{ij}) x_j, \quad (16.95)$$

where δ_{ij} is the Kronecker delta. This can be written in vector form as

$$\frac{d\mathbf{x}}{dt} = \beta \mathbf{M}\mathbf{x}, \quad (16.96)$$

¹⁷This contrasts with the approach we took in Section 16.3.1 where all nodes remained infected for the same amount of time and then recovered. Thus the model studied in this section is not exactly the same as that of Section 16.3.1, being more similar to the traditional SIR model of Section 16.1.2. We will see some minor consequences of this difference shortly.

where \mathbf{M} is the $n \times n$ symmetric matrix

$$\mathbf{M} = \mathbf{A} - \frac{\gamma}{\beta} \mathbf{I}. \quad (16.97)$$

As before we can write \mathbf{x} as a linear combination of eigenvectors, though they are eigenvectors of \mathbf{M} rather than of the adjacency matrix as in the case of the SI model. But now we notice a useful thing: since \mathbf{M} differs from the adjacency matrix only by a multiple of the identity matrix, it has the same eigenvectors \mathbf{v}_r as the adjacency matrix:

$$\mathbf{M}\mathbf{v}_r = \mathbf{A}\mathbf{v}_r - \frac{\gamma}{\beta} \mathbf{I}\mathbf{v}_r = \left(\kappa_r - \frac{\gamma}{\beta} \right) \mathbf{v}_r. \quad (16.98)$$

Only the eigenvalue has been shifted downward by γ/β . Thus the equivalent of Eq. (16.62) is now

$$\mathbf{x}(t) = \sum_{r=1}^n a_r(0) \mathbf{v}_r e^{(\beta\kappa_r - \gamma)t}. \quad (16.99)$$

Note that the exponential constant now depends on $\beta\kappa_r - \gamma$ and so is a function not only of the adjacency matrix and the infection rate but also of the recovery rate, as we would expect—the faster people recover from infection, the slower the disease will spread.

Again the fastest growing term in (16.99) is the one corresponding to the most positive eigenvalue κ_1 of the adjacency matrix and the corresponding eigenvector gives the eigenvector centralities, so an individual's probability of infection at early times is proportional to eigenvector centrality. Note, however, that if γ becomes sufficiently large it is now possible for the exponential constant in the leading term to become negative, meaning that the term decays exponentially rather than grows. Moreover, if the leading term decays then so necessarily do all other terms, and so the total number of infected individuals will decay over time and the disease will die out without causing an epidemic.

The point at which this happens is the epidemic threshold for our model and it occurs at $\beta\kappa_1 - \gamma = 0$, or equivalently

$$\frac{\beta}{\gamma} = \frac{1}{\kappa_1}. \quad (16.100)$$

Thus the position of the epidemic threshold depends on the leading eigenvalue of the adjacency matrix. If the leading eigenvalue is small, then the probability of infection β must be large, or the recovery rate γ small, for the disease to spread. In other words, a small value of κ_1 makes it harder for the disease to spread and a large value easier. This makes intuitive sense, since large values

of κ_1 correspond to denser adjacency matrices and smaller values to sparser ones.

As in the case of the SI model, Eqs. (16.92) to (16.94) are only approximate, because they neglect correlations between the states of adjacent nodes. As before we can allow for these correlations by using a pair approximation, but here we take a different approach and consider instead the equivalent for the SIR model of the degree-based methods of Section 16.5.2.¹⁸

16.6.1 DEGREE-BASED APPROXIMATION FOR THE SIR MODEL

As in our treatment of the SI model in Section 16.5.2, we can develop an approximate solution for the behavior of the SIR model by assuming that all nodes with the same degree behave the same way [342]. Again we consider the example of the configuration model and focus on outbreaks taking place in the giant component of the network. We define $s_k(t)$, $x_k(t)$, and $r_k(t)$ to be the probabilities that a node with degree k is susceptible, infected, or recovered, respectively, at time t . In order for a susceptible node i to become infected it must catch the disease from an infected neighbor j , and for j to be infected it must have caught the disease from one of its neighbors other than i (since i is susceptible). This means, as before, that node j 's probability of being infected is given by x_k , but with k equal to the excess degree of j , which is one less than the total degree. Then, following the arguments of Section 16.5.2, the rate at which the probability of being susceptible decreases is given by the same equation as before, Eq. (16.78):

$$\frac{ds_k}{dt} = -\beta k v s_k, \quad (16.101)$$

where v is the average probability that a neighbor is infected:

$$v(t) = \sum_{k=0}^{\infty} q_k x_k(t). \quad (16.102)$$

¹⁸We can see that the approach of this section cannot be exactly correct from the behavior of Eq. (16.100) on very sparse networks. On a vanishingly sparse network, with only a very few edges and no giant component, κ_1 becomes very small, though still non-zero. On such a network Eq. (16.100) implies that we could, nonetheless, have an epidemic if β is very large or γ very small. Clearly this is nonsense—there can be no epidemic in a network with no giant component. Thus the equation cannot be exactly correct.

And the equations for x_k and r_k are

$$\frac{dx_k}{dt} = \beta k v s_k - \gamma x_k, \quad (16.103)$$

$$\frac{dr_k}{dt} = \gamma x_k. \quad (16.104)$$

We can solve these equations exactly by a combination of the methods of Sections 16.1.2 and 16.5.2. First, we observe that the probability that i 's neighbor j is recovered depends only on the probability that it was previously infected. Since its probability of being infected is x_k (where k is the excess degree), its probability of being recovered is therefore r_k , since by hypothesis all nodes with the same degree behave the same way. Then the average probability that a neighbor is recovered is

$$w(t) = \sum_{k=0}^{\infty} q_k r_k(t). \quad (16.105)$$

Using Eqs. (16.102) and (16.104), we now find that

$$\frac{dw}{dt} = \sum_{k=0}^{\infty} q_k \frac{dr_k}{dt} = \gamma \sum_{k=0}^{\infty} q_k x_k = \gamma v, \quad (16.106)$$

which we use to eliminate v from Eq. (16.101), giving

$$\frac{ds_k}{dt} = -\frac{\beta}{\gamma} k \frac{dw}{dt} s_k. \quad (16.107)$$

This equation can be integrated to give

$$s_k = s_0 \exp\left(-\frac{\beta}{\gamma} k w\right), \quad (16.108)$$

where we have fixed the constant of integration so that at $t = 0$ all nodes have the same probability s_0 of being susceptible and there are no recovered nodes ($w = 0$).

Equation (16.108) implies that s_k is proportional to a power of a universal k -independent function $u(t)$:

$$s_k(t) = s_0 [u(t)]^k, \quad (16.109)$$

where in this case

$$u(t) = e^{-\beta w / \gamma}. \quad (16.110)$$

Then, noting that $s_k + x_k + r_k = 1$, we have

$$\begin{aligned} v(t) &= \sum_k q_k x_k = \sum_k q_k (1 - r_k - s_k) = 1 - w(t) - s_0 \sum_k q_k u^k \\ &= 1 + \frac{\gamma}{\beta} \ln u - s_0 g_1(u), \end{aligned} \tag{16.111}$$

where the value of $w(t)$ comes from rearranging Eq. (16.110) and g_1 is the generating function for the excess degree distribution, as usual.

Substituting (16.111) into (16.106) and again using the value of w from Eq. (16.110), we get

$$\frac{du}{dt} = -\beta u \left[1 + \frac{\gamma}{\beta} \ln u - s_0 g_1(u) \right]. \tag{16.112}$$

This is the equivalent for the SIR model of Eq. (16.84) and differs from that equation only by the new term in $\ln u$ on the right-hand side.

Like (16.84), Eq. (16.112) is a first-order linear differential equation in u and hence can, at least in principle, be solved by direct integration starting from the initial condition $u = 1$ at $t = 0$ (which comes from Eq. (16.110) and the fact that $w = 0$ at $t = 0$ —see Eq. (16.105)). Once we have $u(t)$, the probability s_k of a node being susceptible is given by Eq. (16.109), or we can write the total fraction of susceptibles as

$$s(t) = \sum_k p_k s_k = s_0 \sum_k p_k u^k = s_0 g_0(u). \tag{16.113}$$

Solving for x_k and r_k requires a little more work but with perseverance it can be achieved.¹⁹ Figure 16.14 shows the equivalent of Fig. 16.13 for nodes of a range of degrees. As we can see, the solution has the expected form, with the number of infected individuals rising, peaking, then dropping off as the system evolves to a final state in which some fraction of the population is recovered from the disease and some fraction has never caught it (and never will). Among nodes of different degrees the number infected goes up sharply with degree, as we would expect.

¹⁹We observe that

$$\frac{d}{dt} (e^{\gamma t} x_k) = e^{\gamma t} \left(\frac{dx_k}{dt} + \gamma x_k \right) = e^{\gamma t} \beta k v s_k,$$

where we have used Eq. (16.103) in the second equality. Integrating with respect to t and using Eqs. (16.102) and (16.109), we then have

$$x_k(t) = e^{-\gamma t} \left[x_0 + \beta k s_0 \int_0^t e^{\gamma t'} [u(t')]^k \left[1 + \frac{\gamma}{\beta} \ln u(t') - s_0 g_1(u(t')) \right] dt' \right],$$

and $r_k = 1 - s_k - x_k$.

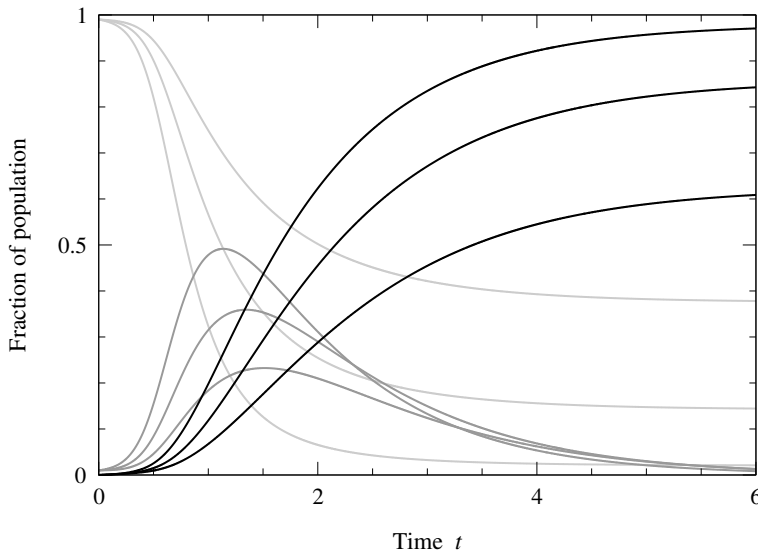


Figure 16.14: Fractions of susceptible, infected, and recovered nodes of various degrees in the SIR model. The fraction of nodes of degree k that are susceptible (light gray), infected (darker gray), and recovered (black) as a function of time for $k = 1, 2, 4$ and $\beta = \gamma = 1$ on a network with an exponential degree distribution $p_k = Ca^k$ with $a = 0.8$. The highest values of k give the fastest growing numbers of infected and recovered nodes and the lowest values the slowest growing.

Depending on the degree distribution (and hence the form of $g_1(u)$) it may not be possible to solve Eq. (16.112) in closed form, but even when it isn't the equations can still shed light on features of the epidemic. Consider, for example, the long-time behavior. In the limit of long time we expect that the number of infected individuals will go to zero, leaving some individuals recovered and some who have never caught the disease. At $t = \infty$ the total fraction $r(t)$ of recovered individuals is equal to the fraction who have ever had the disease, i.e., the overall size of the outbreak, and is given by

$$r(\infty) = 1 - s(\infty) = 1 - g_0(u(\infty)), \quad (16.114)$$

where we have set $s_0 = 1$ as before on the assumption that the system is large and the number of initially infected individuals small.

We can find the limiting value of u by setting $du/dt = 0$ in Eq. (16.112) to give

$$1 + \frac{\gamma}{\beta} \ln u - g_1(u) = 0. \quad (16.115)$$

In the special case where the outbreak is small, so that the final value of u is close to 1, we can expand $\ln u = \ln[1 + (u - 1)] \approx u - 1$ and Eq. (16.115) becomes

$$u \approx 1 - \frac{\beta}{\gamma} + \frac{\beta}{\gamma} g_1(u). \quad (16.116)$$

Equations (16.114) and (16.116) are similar in form to Eqs. (16.29) and (16.30), which give the final size of the outbreak in our treatment of the SIR model using percolation theory. The reason why Eq. (16.116) is only approximate in the present case where Eq. (16.30) was exact is that the model treated in this section is slightly different from the one treated earlier, having (as discussed previously) a constant probability γ per unit time of recovery from disease for each infected individual as opposed to a fixed infection time (see footnote 17 on page 660).

We can also examine the early-time behavior of the outbreak. Recall that $u = 1$ at $t = 0$ so at early times the value of u will be close to 1 (and slightly smaller, since u is a probability and cannot be greater than 1). Setting $u = 1 - \epsilon$ in Eq. (16.112) and keeping terms to leading order in ϵ and x_0 we get

$$\frac{d\epsilon}{dt} = \beta x_0 + [\beta g'_1(1) - \gamma] \epsilon, \quad (16.117)$$

which means that

$$u(t) = 1 - \epsilon(t) = 1 - \frac{\beta x_0}{\beta g'_1(1) - \gamma} [e^{(\beta g'_1(1) - \gamma)t} - 1]. \quad (16.118)$$

This is similar to Eq. (16.90) for the SI model, except for the inclusion of γ . The fraction of susceptible degree- k nodes at early times is given by

$$s_k(t) = u^k = (1 - \epsilon)^k \approx 1 - k\epsilon, \quad (16.119)$$

and the number of cases of the disease, infected and recovered, which is proportional to $1 - s_k$, grows exponentially as $ke^{[\beta g'_1(1) - \gamma]t}$.

The epidemic threshold for the model is the line that separates an initially growing number of cases of the disease from an initially decreasing one and is given in this case by the point at which the exponential constant in Eq. (16.118) equals zero, which gives

$$\frac{\beta}{\gamma} = \frac{1}{g'_1(1)}. \quad (16.120)$$

This result is similar in form to Eq. (16.100) for the epidemic threshold on a general network,²⁰ but with the leading eigenvalue of the adjacency matrix κ_1

²⁰And like Eq. (16.100) it is also clearly wrong on sparse networks for the same reasons—see footnote 18 on page 662.

replaced with $g'_1(1)$. It also looks similar to Eq. (16.31) for the percolation threshold for bond percolation, but this similarity is somewhat deceptive. In fact, the result most nearly corresponding to this one in the percolation treatment is Eq. (16.32). If we equate our recovery rate γ with the reciprocal of the infectiousness time τ in that previous treatment, then the two are roughly equivalent when the epidemic threshold is low, meaning either that β is small or that γ is large. If the threshold is higher then the match between the two models is poorer, which is again a result of the fact that the models are defined in slightly different ways.

16.7 TIME-DEPENDENT PROPERTIES OF THE SIS MODEL

So far in this chapter we have said little about the behavior of the SIS model on a network, or any other model that incorporates reinfection. As mentioned in Section 16.4, the equivalence to percolation processes that we used to analyze the SIR model does not apply when reinfection is involved. On the other hand, the differential equation methods of the preceding sections can be extended to the SIS model in a straightforward fashion. By analogy with Eqs. (16.92) to (16.94) we have

$$\frac{ds_i}{dt} = -\beta s_i \sum_j A_{ij} x_j + \gamma x_i, \quad (16.121a)$$

$$\frac{dx_i}{dt} = \beta s_i \sum_j A_{ij} x_j - \gamma x_i \quad (16.121b)$$

for the SIS model. Caveats similar to those for previous models apply here: these equations ignore correlations between the states of adjacent nodes and hence are only an approximation.

Equations (16.121a) and (16.121b) are not independent since $s_i + x_i = 1$, so only one of them is needed to form a solution. Taking the second and eliminating s_i we get

$$\frac{dx_i}{dt} = \beta(1 - x_i) \sum_j A_{ij} x_j - \gamma x_i. \quad (16.122)$$

At early times, assuming as before that $x_i(0) = x_0 = c/n$ for all i and constant c , we can drop terms at quadratic order in small quantities to get

$$\frac{dx_i}{dt} = \beta \sum_j A_{ij} x_j - \gamma x_i, \quad (16.123)$$

which is identical to Eq. (16.95) for the SIR model at early times. Hence we can immediately conclude that the early-time behavior of the model is the same (within this approximation), with initially exponential growth and an epidemic threshold given by

$$\frac{\beta}{\gamma} = \frac{1}{\kappa_1}. \quad (16.124)$$

(See Eq. (16.100).) Also as in the SIR model the probability of infection of a given node at early times will be proportional to the node's eigenvector centrality.

At late times we expect the probability of infection to settle to a constant endemic level, which we can calculate by setting $dx_i/dt = 0$ in Eq. (16.122) and rearranging, to give

$$x_i = \frac{\sum_j A_{ij}x_j}{\gamma/\beta + \sum_j A_{ij}x_j}. \quad (16.125)$$

Typically we cannot derive a closed-form solution for x_i from this expression, but we can solve it numerically by iteration starting from a random initial guess. We can also see the general form the solution will take by considering limiting cases. If β/γ is large, meaning that we are well above the epidemic threshold given in Eq. (16.124), then we can ignore the term γ/β in the denominator and $x_i \approx 1$ for all i , meaning that essentially all nodes will be infected all the time. This makes good sense since if β/γ is large then the rate of infection is very high while the rate of recovery is small.

Conversely, if β/γ is only just above the epidemic threshold level set by Eq. (16.124) then x_i will be small—the disease only just manages to stay alive—and we can ignore the sum in the denominator of Eq. (16.125) so that

$$x_i \approx \frac{\beta}{\gamma} \sum_j A_{ij}x_j, \quad (16.126)$$

or

$$\kappa_1 x_i \approx \sum_j A_{ij}x_j, \quad (16.127)$$

where we have used Eq. (16.124). This implies that x_i is proportional to the leading eigenvector of the adjacency matrix or, equivalently, proportional to the eigenvector centrality. (Note that this is at late times so this result is distinct from our earlier finding that x_i is proportional to eigenvector centrality at early times.)

Thus, the long-time endemic disease behavior of the SIS model varies from a regime just above the epidemic threshold in which the probability of a node being infected is proportional to its eigenvector centrality, to a regime well above the threshold in which essentially every node is infected at all times.

16.7.1 DEGREE-BASED APPROXIMATION FOR THE SIS MODEL

We can also write down approximate equations for the evolution of the SIS model in which, as in Sections 16.5.2 and 16.6.1, we assume that the probability of infection is the same for all nodes with a given degree. Focusing once again on configuration model networks, the equivalent of Eqs. (16.101) and (16.103) is

$$\frac{ds_k}{dt} = -\beta k v s_k + \gamma x_k, \quad (16.128a)$$

$$\frac{dx_k}{dt} = \beta k v s_k - \gamma x_k, \quad (16.128b)$$

where the variables s_k and x_k are as before, and v is the average probability that a neighbor is infected:

$$v(t) = \sum_{k=0}^{\infty} q_k x_k(t). \quad (16.129)$$

As before Eqs. (16.128a) and (16.128b) are not independent and only one is needed to form a solution. Taking the second and rewriting it using $s_k = 1 - x_k$ we get

$$\frac{dx_k}{dt} = \beta k v (1 - x_k) - \gamma x_k. \quad (16.130)$$

Unfortunately, there is no known complete solution to this equation but we can once again find its behavior at early and late times.

Assuming, as previously, that our epidemic starts off with only a single case or a small number of cases, the probability x_k of being infected at early times is c/n for constant c and hence small in the limit of large n . Dropping terms of second order in small quantities then gives us the linearized equation

$$\frac{dx_k}{dt} = \beta k v - \gamma x_k, \quad (16.131)$$

which can be rewritten using an integrating factor to read

$$\frac{d}{dt} (e^{\gamma t} x_k) = e^{\gamma t} \beta k v + \gamma e^{\gamma t} x_k = \beta k e^{\gamma t} v, \quad (16.132)$$

and hence integrated to give

$$x_k(t) = \beta k \int_0^t e^{\gamma(t-t')} v(t') dt'. \quad (16.133)$$

Thus $x_k(t)$ for short times takes the form

$$x_k = k u(t), \quad (16.134)$$

where $u(t)$ is a universal, k -independent function. Substituting into Eqs. (16.129) and (16.131), we then have

$$v(t) = u(t) \sum_{k=0}^{\infty} kq_k = g_1'(1)u(t), \quad (16.135)$$

and

$$\frac{du}{dt} = [\beta g_1'(1) - \gamma]u(t). \quad (16.136)$$

Thus we have exponential growth or decay of the epidemic at early times, with the epidemic threshold that separates the two falling at the point where $\beta g_1'(1) - \gamma = 0$, or

$$\frac{\beta}{\gamma} = \frac{1}{g_1'(1)}, \quad (16.137)$$

just as for the SIR model (see Eq. (16.120)).

At late times the disease settles into an endemic state with a constant fraction of the population infected. We can solve for this state by setting $dx_k/dt = 0$ for all k in Eq. (16.130) to give

$$x_k = \frac{kv}{kv + \gamma/\beta}. \quad (16.138)$$

Substituting this expression into Eq. (16.129), we then find that

$$\sum_{k=0}^{\infty} \frac{kq_k}{kv + \gamma/\beta} = 1. \quad (16.139)$$

If we can solve this equation for v then we can then get x_k from Eq. (16.138).

Unfortunately, there is, in general, no closed-form solution for Eq. (16.139), although it can be solved numerically for any given q_k . What we can say is that, given the degree distribution, v at late times is a function solely of β/γ (or γ/β if you prefer) and hence x_k is solely a function of β/γ and k . Moreover, in order for Eq. (16.139) to be satisfied v must be an increasing function of β/γ —as β gets larger or γ smaller, v must increase in order to keep the sum in the equation equal to 1. This means that x_k will also be an increasing function of β/γ . (Equation (16.138) implies that it is an increasing function of k as well.) Thus, the equations give us a qualitative picture of the behavior of the SIS model, although quantitative details require a numerical solution.

We have in this chapter only brushed the surface of what is possible in the modeling of epidemics spreading across networks. We can extend our calculations to more complicated network structures, such as networks with degree

correlations, networks with transitivity, networks with community structure, and even epidemics on empirically observed networks. More complicated models of the spread of infection are also possible, such as the SIRS model mentioned in Section 16.1.6, as well as models that incorporate birth, death, or geographic movement of individuals [21, 232]. In recent years, researchers have developed extremely sophisticated computer models of disease spread using complex simulations of the behavior patterns of human populations, including models of entire cities down to the level of individual people, cars, and buildings [165], and models of the international spread of disease that incorporate detailed data on the flight patterns and timetables of international airlines [118].

EXERCISES

16.1 Consider the bond percolation model of an epidemic in Section 16.3.1 and suppose that for a particular value of ϕ the giant cluster occupies a fraction S of the network. What is the probability of an epidemic outbreak if the disease starts simultaneously at c different nodes, chosen uniformly and independently at random from the whole network? Note that this probability tends exponentially to 1 as c gets larger. The chances of avoiding an epidemic become slim when a disease starts at many points simultaneously.

16.2 Footnote 6 on page 613 discusses a variant of the SIR model in which people die of the disease rather than recovering. Assuming, as in the ordinary model, that an individual makes contacts with others at a rate β per unit time, but that those contacts are with living people only, write down a modified version of the SIR equations (16.9) for this variant of the model.

16.3 Consider an epidemic SIR outbreak (i.e., an outbreak that starts in the giant cluster of the corresponding percolation process) on a configuration model network with exponential degree distribution $p_k = (1 - a)a^k$ with $a < 1$. You can assume that the network is large.

- a) Using the results of Section 15.2.1 write down an expression for the probability u appearing in Eq. (16.29) in terms of ϕ and a .
- b) Hence find an expression for the probability that a node is infected by the disease if it has degree k .
- c) Evaluate this probability for the case $a = 0.4$ and $\phi = 0.9$, for $k = 0, 1$, and 10.

16.4 Consider the SIR model on a configuration model network where all nodes have degree four (also known as a random 4-regular graph).

- a) What is the critical value ϕ_c of the transmission probability at the epidemic threshold?
- b) When the transmission probability is $\phi = \frac{1}{2}$, show that if an epidemic happens, the fraction S of nodes infected is

$$S = \frac{3\sqrt{5} - 5}{2}.$$

16.5 Create a computer simulation, in the programming language of your choice, of the spread of an SIR-type disease over a network. Your program should perform the following steps:

- a) Generate a Poisson random graph of the type $G(n, m)$ with $n = 10\,000$ nodes and $m = 25\,000$ edges, and select one node at random to be the single initial carrier of the disease.
- b) On each time step, every currently infected node spreads the disease to each of its currently susceptible neighbors with independent probability $\phi = 0.4$, then recovers and remains in the recovered state indefinitely thereafter. (This particular variant of SIR, in which nodes remain infective for exactly one time step, is called the Reed–Frost model.)

Run your program for many time steps until no infected nodes remain in the network and make a graph showing, on the same axes, the number of susceptible, infected, and recovered nodes as a function of time. (You may need to run the program several times to find a good example. Sometimes you will find that the disease dies out after only infecting a few individuals and no epidemic occurs.)

16.6 Consider the spread of an SIR-type disease on a configuration model network in which some fraction of the individuals have been vaccinated against the disease. We can model this situation using a joint site/bond percolation model in which a fraction ϕ_s of the nodes are occupied, to represent those not vaccinated, and a fraction ϕ_b of the edges are occupied to represent those along which contact takes place.

- a) Show that the fraction S of individuals infected in the limit of long time is given by the solution of the equations

$$S = \phi_s [1 - g_0(u)], \quad u = 1 - \phi_s \phi_b + \phi_s \phi_b g_1(u),$$

where g_0 and g_1 are the generating functions for the degree distribution and excess degree distribution, as usual.

- b) Show that for a given probability of contact ϕ_b the minimum fraction of individuals that need to be vaccinated to prevent the disease from spreading is $1 - 1/[\phi_b g_1'(1)]$.

16.7 We have been concerned in this chapter primarily with *epidemic* disease outbreaks, meaning outbreaks that affect a finite fraction of all individuals in a network. Consider, by contrast, a small SIR outbreak—an outbreak that corresponds to one of the non-giant percolation clusters in the bond percolation approach of Section 16.3.1—occurring on a configuration model network with degree distribution p_k .

- a) What is the probability of such an outbreak occurring if the disease starts at a node chosen uniformly at random from the whole network (including nodes both within and outside the giant component)?
- b) Show that if the probability of transmission along an edge is ϕ , then the generating function $h_0(z)$ for the probability π_s that the outbreak has size s is given by the equations

$$h_0(z) = z g_0(h_1(z)), \quad h_1(z) = 1 - \phi + \phi z g_1(h_1(z)),$$

where $g_0(z)$ and $g_1(z)$ are the generating functions for the degree distribution and excess degree distribution, respectively.

- c) What is the mean size of such an outbreak?

16.8 Consider an SI-type epidemic spreading on the giant component of a k -regular random graph, i.e., a configuration model network in which all nodes have the same degree k . Assume that some number c of nodes, chosen at random, are infected at time $t = 0$.

- a) Show using the results of Section 16.5 that the probability of infection of every node increases at short times as $e^{\beta k t}$.
- b) Show that within the first-order moment closure approximation of Eq. (16.55) the average probability of infection x of every node is the same and give the differential equation it satisfies.
- c) Hence show that

$$x(t) = \frac{c e^{\beta k t}}{n - c + c e^{\beta k t}}.$$

- d) Find the time at which the “inflection point” of the epidemic occurs, the point at which the rate of appearance of new disease cases stops increasing and starts decreasing.

16.9 Consider a configuration model network containing nodes of degrees 1, 2, and 3 only, such that the fractions of nodes of each degree in the giant component are $p_1 = 0.3$, $p_2 = 0.3$, and $p_3 = 0.4$.

- a) Find an expression for the generating function $g_1(z)$ of the excess degree distribution.
- b) Hence, by solving Eq. (16.84), find an expression for t as a function of u for an SI epidemic on the giant component of the network, assuming that $s_0 \approx 1$, and with initial condition $u(0) = 1 - \epsilon$, where ϵ is small.
- c) Show that in the limit of long times the number of susceptibles falls off in proportion to $e^{-21\beta t/2}$.

16.10 Consider a configuration model network with any degree distribution.

- a) By comparing Eqs. (16.100) and (16.120), and using the definition of the generating function g_1 , give an expression for the largest eigenvalue κ_1 of the adjacency matrix on a configuration model network in terms of the mean and mean-square degree in the network.

- b) In Section 17.4 it is shown that the largest eigenvalue of the adjacency matrix of any network satisfies $\kappa_1 \geq \sqrt{\langle k^2 \rangle}$ (see Eq. (17.89)). For a network with a Poisson degree distribution, show that your expression for κ_1 violates this inequality. How do you explain this?

CHAPTER 17

DYNAMICAL SYSTEMS ON NETWORKS

A discussion of dynamical systems on networks and methods for their analysis

THE epidemic models of Chapter 16 are one example of the more general concept of dynamical systems on networks. A dynamical system is any system whose state, as represented by some set of variables, changes over time according to some given rules or equations. Dynamical systems come in continuous- and discrete-time varieties and can be either deterministic or stochastic. The differential-equation-based models we looked at for epidemics, for instance, are continuous-time dynamical systems because the equations describe the continuous-time evolution of their variables. They are also deterministic because the equations exactly determine the values of all variables for all time: there is no random or external element in the equations whose value is not known in advance. On the other hand, an explicit computer simulation of, say, an SI epidemic on a network would be a stochastic dynamical system and might use either continuous or discrete time. The stochastic element in this case is the random infection of susceptible individuals by their infectious neighbors. And time might be represented in discrete time-steps, such as days, or it might be continuous, depending on the decision of the researcher.

Many other real-world processes—or models of real-world processes—can be thought of as dynamical systems on networks. The spread of news or information between friends, the movement of money through an economy, the flow of traffic on roads, data over the Internet, or electricity over the grid, the evolution of populations in an ecosystem, the changing concentrations of metabolites

in a cell, and many other systems of scientific interest can be represented as dynamical systems of one kind or another evolving on an appropriate network.

In other, non-network contexts, the theory of dynamical systems is a well-developed branch of mathematics and physics—see, for example, the book by Strogatz [441]. In this chapter we delve into some of this theory and show how it can be applied to dynamical systems on networks. Necessarily our introduction covers only a small part of what could be said; dynamical systems is a topic of entire books in its own right [45,390,441]. But the material covered here gives a flavor of the kinds of calculation that are possible.

17.1 DYNAMICAL SYSTEMS

Our discussion in this chapter will concentrate primarily on deterministic dynamical systems of continuous real-valued variables evolving in continuous time t . We begin by introducing some of the basic ideas in a non-network context, then we extend these ideas to networks.

A simple (non-network) example of a continuous dynamical system is a system described by a single real variable $x(t)$ that evolves according to a first-order differential equation

$$\frac{dx}{dt} = f(x), \quad (17.1)$$

where $f(x)$ is some specified function of x . Typically we will also give an initial condition that specifies the value x_0 taken by x at some initial time t_0 .

The fully-mixed SI model of Section 16.1.1 is an example of a dynamical system of this kind, having a single variable $x(t)$ representing the fraction of infected individuals in the system at time t and obeying the equation

$$\frac{dx}{dt} = \beta x(1 - x). \quad (17.2)$$

(See Eq. (16.5).) Thus in this case we have $f(x) = \beta x(1 - x)$.

One can also have dynamical systems of two variables:

$$\frac{dx}{dt} = f(x, y), \quad \frac{dy}{dt} = g(x, y), \quad (17.3)$$

and the approach can be extended to larger numbers of variables as well. When we come to consider systems on networks we will put separate variables on each node of the network.

One could also imagine making the functions on the right-hand side of our equations depend explicitly on time t :

$$\frac{dx}{dt} = f(x, t). \quad (17.4)$$

This, however, can be regarded as merely a special case of Eq. (17.3). If we write

$$\frac{dx}{dt} = f(x, y), \quad \frac{dy}{dt} = 1, \quad (17.5)$$

with initial condition $y(0) = 0$, then we have $y = t$ for all times and $dx/dt = f(x, t)$ as required. By this trick it is always possible to turn equations with explicit dependence on t into equations without explicit dependence on t but with one extra variable. For this reason we will confine ourselves in this chapter to systems with no explicit dependence on t .

Another possible generalization would be to consider systems governed by equations containing higher derivatives, such as second derivatives, or functions of derivatives. But these can also be reduced to simpler cases by introducing extra variables. For instance, the equation

$$\frac{d^2x}{dt^2} + \left(\frac{dx}{dt}\right)^2 - \frac{dx}{dt} = f(x) \quad (17.6)$$

can be transformed by introducing a new variable $y = dx/dt$ so that we have

$$\frac{dx}{dt} = y, \quad \frac{dy}{dt} = f(x) - y^2 + y, \quad (17.7)$$

which is a special case of Eq. (17.3) again.

Thus the study of systems of equations like (17.1) and (17.3) covers a broad range of situations of scientific interest. Let us look at some of the techniques used to analyze such equations.

17.1.1 FIXED POINTS AND LINEARIZATION

Equation (17.1), which involves only the one variable x , can, at least in principle, always be solved by rearranging and integrating to give

$$\int_{x_0}^x \frac{dx'}{f(x')} = t - t_0, \quad (17.8)$$

although in practice the integral may not be known in closed form. For cases with two or more variables it is in general not possible to find a solution. And for the network examples that we will be studying shortly the number of variables is typically very large, so that, unless we are lucky (as we were with some of the epidemiological models of the previous chapter), full analytic solutions are unlikely to be forthcoming.

One alternative is to solve the equations numerically, on a computer, and this can often give useful insight, but let us not give up on analytic approaches just

yet. There is in fact a well-developed set of techniques for understanding how dynamical systems work without first solving their equations exactly. Most of those techniques focus on the properties of fixed points.

A *fixed point* is a steady state of the system—any value of the variable or variables for which the system is stationary and doesn't change over time. In the one-variable system, Eq. (17.1), for example, a fixed point is any point $x = x^*$ for which the function on the right-hand side of the equation is zero,

$$f(x^*) = 0, \quad (17.9)$$

so that $dx/dt = 0$ and x doesn't move. If, in the evolution of the system, x ever reaches a fixed point, then it will remain there forever. The fixed points of a one-variable system can be found simply by solving $f(x) = 0$ for x .

In a two-variable system like Eq. (17.3) a fixed point is a pair of values (x^*, y^*) such that $f(x^*, y^*) = 0$ and $g(x^*, y^*) = 0$, making $dx/dt = dy/dt = 0$ so that both variables stand still at this point. Thus we can find the fixed points by solving the simultaneous equations $f(x, y) = 0$ and $g(x, y) = 0$. In general, the fixed points of an n -variable system are found by solving n simultaneous equations.

As an example, consider the SI model of Eq. (17.2). Putting $f(x) = 0$ in this model gives us $\beta x(1 - x) = 0$, which has solutions $x = 1$ and $x = 0$ for the fixed points. We can see immediately what these fixed points mean in epidemiological terms. The first at $x = 1$ represents the steady state in which everyone in the system is infected. Clearly, once everyone is infected the system doesn't change any more, because there is no one else to infect and because in the SI model no one recovers either. The second fixed point $x = 0$ corresponds to the state of the system where no one is infected. In this state no one will ever become infected, since there is no one to catch the disease from, so again we have a steady state.

The importance of fixed points in the study of dynamical systems derives from two key features of these points: first, they are relatively easy to find, and second, it is straightforward to determine the dynamics of the system when it is close to, but not exactly at, a fixed point. The dynamics close to a fixed point is found by expanding about the point as follows.

Consider first a simple one-variable system obeying Eq. (17.1). We represent the value of x close to a fixed point at x^* by writing $x = x^* + \epsilon$ where ϵ , which represents our distance from the fixed point, is small. Then

$$\frac{dx}{dt} = \frac{d\epsilon}{dt} = f(x^* + \epsilon). \quad (17.10)$$

Now we perform a Taylor expansion of the right-hand side about the point

$x = x^*$ to get

$$\frac{d\epsilon}{dt} = f(x^*) + \epsilon f'(x^*) + O(\epsilon^2), \quad (17.11)$$

where f' represents the derivative of f with respect to its argument. Neglecting terms of order ϵ^2 and smaller and noting that $f(x^*) = 0$ by definition (see Eq. (17.9)), we then have

$$\frac{d\epsilon}{dt} = \epsilon f'(x^*). \quad (17.12)$$

This is a linear first-order differential equation with solution

$$\epsilon(t) = \epsilon(0)e^{\lambda t}, \quad (17.13)$$

where

$$\lambda = f'(x^*). \quad (17.14)$$

The quantity λ is called a *Lyapunov exponent*. Note that it is just a simple number, which we can calculate provided we know the position x^* of the fixed point and the function $f(x)$. Depending on the sign of λ , Eq. (17.13) tells us that our distance ϵ from the fixed point will either grow or decay exponentially in time. Thus this analysis allows us to classify our fixed points into two types. An *attracting fixed point* is one with $\lambda < 0$, so that points close by are attracted towards the fixed point and eventually flow into it. A *repelling fixed point* is one with $\lambda > 0$, so that points close by are repelled away. In between these two types there is a special case when $\lambda = 0$ exactly. Fixed points with $\lambda = 0$ are usually still either attracting or repelling,¹ but one cannot tell which from the analysis here; one must retain some of the higher-order terms that we dropped in Eq. (17.11) to determine what happens.

Analysis of the kind represented by Eq. (17.12) is known as *linear stability analysis*. It can be applied to systems with two or more variables as well. Consider, for instance, a dynamical system governed by equations of the form of Eq. (17.3), with a fixed point at (x^*, y^*) , meaning that

$$f(x^*, y^*) = 0, \quad g(x^*, y^*) = 0. \quad (17.15)$$

We represent a point close to the fixed point in the two-dimensional x, y space by $x = x^* + \epsilon_x$ and $y = y^* + \epsilon_y$, where ϵ_x and ϵ_y are both assumed small.

¹There are also a couple of other rarer possibilities. A fixed point with $\lambda = 0$ can be *neutral*, meaning it neither attracts nor repels. For example, the choice $f(x) = 0$ for all x has a neutral fixed point at every value of x . Another less trivial possibility is that a fixed point with $\lambda = 0$ may be of *mixed* type, meaning that it attracts on one side and repels on the other. An example is $f(x) = x^2$, which has a fixed point at $x = 0$ that is attracting for $x < 0$ and repelling for $x > 0$.

As before we expand about the fixed point, performing now a double Taylor expansion:

$$\begin{aligned}\frac{dx}{dt} &= \frac{d\epsilon_x}{dt} = f(x^* + \epsilon_x, y^* + \epsilon_y) \\ &= f(x^*, y^*) + \epsilon_x f^{(x)}(x^*, y^*) + \epsilon_y f^{(y)}(x^*, y^*) + \dots,\end{aligned}\quad (17.16)$$

where $f^{(x)}$ and $f^{(y)}$ indicate the derivatives of f with respect to x and y . Making use of Eq. (17.15) and neglecting all higher-order terms in the expansion, we can simplify this expression to

$$\frac{d\epsilon_x}{dt} = \epsilon_x f^{(x)}(x^*, y^*) + \epsilon_y f^{(y)}(x^*, y^*),\quad (17.17)$$

and similarly

$$\frac{d\epsilon_y}{dt} = \epsilon_x g^{(x)}(x^*, y^*) + \epsilon_y g^{(y)}(x^*, y^*).\quad (17.18)$$

We can combine these two equations and write them in matrix form as

$$\frac{d\epsilon}{dt} = \mathbf{J}\epsilon,\quad (17.19)$$

where ϵ is the two-component vector (ϵ_x, ϵ_y) and \mathbf{J} is the Jacobian matrix

$$\mathbf{J} = \begin{pmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{pmatrix},\quad (17.20)$$

where the derivatives are all evaluated at the fixed point. For systems of three or more variables we can employ the same approach and again we arrive at Eq. (17.19), but with the matrices and vectors being larger, having as many rows and columns as there are variables.

Equation (17.19) is again a linear first-order differential equation but its solution is more complicated than for the one-variable equivalent. The simplest case occurs when the Jacobian matrix is diagonal, in which case the equation can be written

$$\begin{pmatrix} \frac{d\epsilon_x}{dt} \\ \frac{d\epsilon_y}{dt} \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} \epsilon_x \\ \epsilon_y \end{pmatrix},\quad (17.21)$$

where λ_1 and λ_2 are real numbers. Then the variables ϵ_x and ϵ_y separate from one another thus:

$$\frac{d\epsilon_x}{dt} = \lambda_1 \epsilon_x, \quad \frac{d\epsilon_y}{dt} = \lambda_2 \epsilon_y,\quad (17.22)$$

and we can solve for each variable separately to get

$$\epsilon_x(t) = \epsilon_x(0) e^{\lambda_1 t}, \quad \epsilon_y(t) = \epsilon_y(0) e^{\lambda_2 t}, \quad (17.23)$$

or equivalently

$$x(t) = x^* + \epsilon_x(0) e^{\lambda_1 t}, \quad y(t) = y^* + \epsilon_y(0) e^{\lambda_2 t}. \quad (17.24)$$

Thus x and y are independently either attracted or repelled from the fixed point over time, depending on the signs of the two quantities

$$\lambda_1 = \left(\frac{\partial f}{\partial x} \right)_{\substack{x=x^* \\ y=y^*}}, \quad \lambda_2 = \left(\frac{\partial g}{\partial y} \right)_{\substack{x=x^* \\ y=y^*}}. \quad (17.25)$$

These results give rise to a variety of possible behaviors of the system near the fixed point, as shown in Fig. 17.1. If λ_1 and λ_2 are both negative, for instance, then the fixed point will be attracting, while if they are both positive it will be repelling. If they are of opposite signs then we have a new type of point called a *saddle point* that attracts along one axis and repels along the other. In some respects a saddle point is perhaps best thought of as a form of repelling fixed point, since a system that starts near a saddle point will not stay near it, the dynamics being repelled along the unstable direction.

Unless we are very lucky, however, the Jacobian matrix is unlikely to be diagonal. In the general case it will have off-diagonal as well as diagonal elements and the solution method above will not work. With a little more effort, however, we can make progress in this case too. The trick is to find combinations of the variables x and y that move independently, as x and y alone do above.

Consider the linear combinations of variables

$$\xi_1 = a\epsilon_x + b\epsilon_y, \quad \xi_2 = c\epsilon_x + d\epsilon_y. \quad (17.26)$$

In matrix form we can write these as

$$\begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} \epsilon_x \\ \epsilon_y \end{pmatrix}, \quad (17.27)$$

or simply

$$\xi = \mathbf{Q}\epsilon, \quad (17.28)$$

where \mathbf{Q} is the matrix of the coefficients a, b, c, d .

The time evolution of ξ close to the fixed point is given by

$$\frac{d\xi}{dt} = \mathbf{Q} \frac{d\epsilon}{dt} = \mathbf{Q}\mathbf{J}\epsilon = \mathbf{Q}\mathbf{J}\mathbf{Q}^{-1}\xi, \quad (17.29)$$

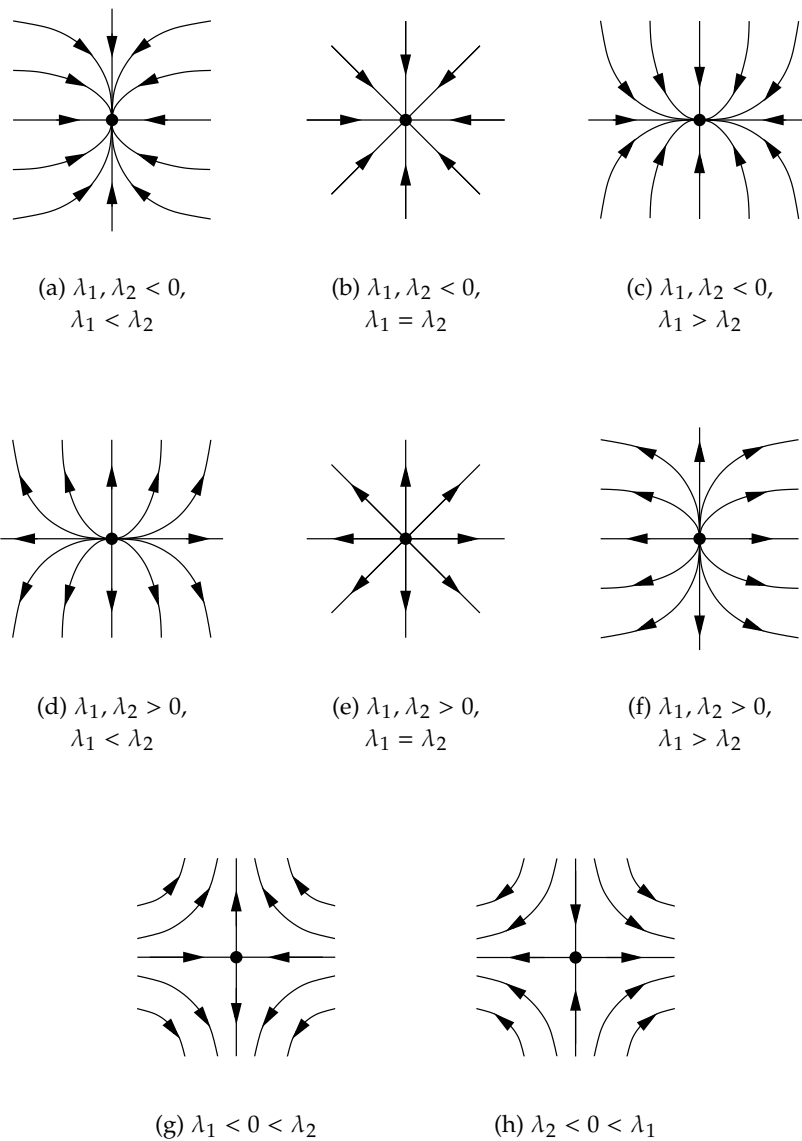


Figure 17.1: Flows in the vicinity of different types of fixed points. The flows around a fixed point in a two-variable dynamical system with a diagonal Jacobian matrix, as described in the text, can take a variety of different forms as shown. (a), (b), and (c) are all attracting fixed points, (d), (e), and (f) are repelling, and (g) and (h) are saddle points.

where we have used Eqs. (17.19) and (17.28). If ξ_1 and ξ_2 are to evolve independently, then we require that the matrix $\mathbf{Q}\mathbf{J}\mathbf{Q}^{-1}$ be diagonal, just as \mathbf{J} itself was in the simple case we studied earlier. Linear algebra then tells us that \mathbf{Q} must be the matrix of eigenvectors of \mathbf{J} . Specifically, since \mathbf{J} is in general asymmetric, \mathbf{Q} is the matrix whose rows are the left eigenvectors of \mathbf{J} and the inverse \mathbf{Q}^{-1} is the matrix whose columns are the right eigenvectors of \mathbf{J} (since the left and right eigenvectors of a matrix are mutually orthogonal). Then Eq. (17.28) (or equivalently Eq. (17.26)) tells us the combinations ξ_1 and ξ_2 that move independently of one another near the fixed point.

These combinations satisfy the equations

$$\frac{d\xi_1}{dt} = \lambda_1 \xi_1, \quad \frac{d\xi_2}{dt} = \lambda_2 \xi_2, \quad (17.30)$$

where λ_1 and λ_2 are the elements of our diagonal matrix, which are also the eigenvalues of \mathbf{J} corresponding to the two eigenvectors. Equation (17.30) has the obvious solution

$$\xi_1(t) = \xi_1(0) e^{\lambda_1 t}, \quad \xi_2(t) = \xi_2(0) e^{\lambda_2 t}. \quad (17.31)$$

The lines $\xi_1 = 0$ and $\xi_2 = 0$ play the role of the axes in Fig. 17.1—they are lines along which we move either directly away from or directly towards the fixed point—and Eq. (17.31) indicates that our distance from the fixed point along these lines will either grow or decay exponentially according to the signs of the two eigenvalues. Since the eigenvectors of an asymmetric matrix are not in general orthogonal to one another, these lines are not in general at right angles, so the flows around the fixed point will look similar to those of Fig. 17.1 but squashed, as shown in Fig. 17.2. Nonetheless, we can still classify our fixed points as attracting, repelling, or saddle points as shown in the figure. Similar analyses can be performed for systems with larger numbers of variables and the basic results are the same: by finding the eigenvectors of the Jacobian matrix we can determine the combinations of variables that move independently and hence solve for the evolution of the system in the vicinity of the fixed point.

There is, however, an additional subtlety that arises for systems of two or more variables that is not found in the one-variable case. The eigenvalues of an asymmetric matrix need not be real. Even if the elements of the matrix itself are real, the eigenvalues can be imaginary or complex. What does it mean if the eigenvalues of the Jacobian matrix in our derivation are complex? Putting such eigenvalues into Eq. (17.31) gives us a solution that *oscillates* around the fixed point, rather than simply growing or decaying. Indeed, the substitution actually gives us a value for ξ_1 and ξ_2 that itself is complex, which looks like it

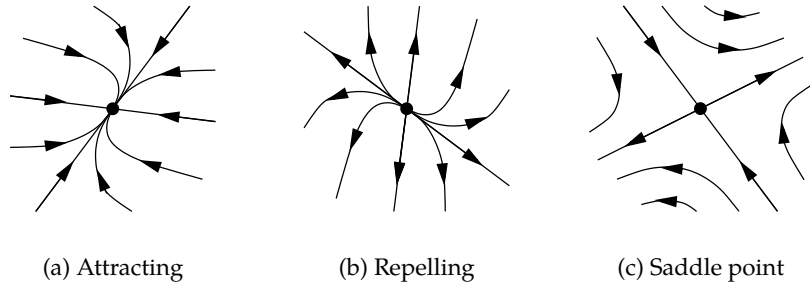


Figure 17.2: Examples of flows around general fixed points. When the Jacobian matrix is not diagonal the flows around a fixed point look like squashed or stretched versions of those in Fig. 17.1.

might be a problem, since the coordinates are supposed to be real. However, our equations are linear, so the real part of that solution is also a solution.

If $\lambda_1 = \alpha + i\omega$, for example, where α and ω are real numbers, then the general real solution for ξ_1 is

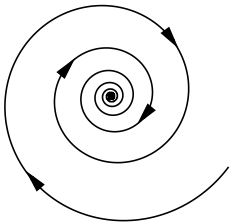
$$\xi_1(t) = \text{Re}[C e^{(\alpha+i\omega)t}] = e^{\alpha t}(A \cos \omega t + B \sin \omega t), \quad (17.32)$$

where A and B are real constants and C is a complex constant. Thus the solution is the product of a part that oscillates and a part that either grows or decays exponentially. For the case of two variables, it turns out that the eigenvalues are always either both real or both complex, and if both are complex then they are complex conjugates of one another. In the latter case, both ξ_1 and ξ_2 then have this combined behavior of oscillation with exponential growth or decay, with the same frequency ω of oscillation and the same rate of growth or decay. The net result is a trajectory that describes a spiral around the fixed point. Depending on whether α is positive or negative, the spiral either goes outward from the fixed point or inward. If it goes inward, i.e., if $\alpha < 0$, then the fixed point is a stable one; otherwise, if $\alpha > 0$, it is unstable. Thus stability is in this case determined solely by the real part of the eigenvalues. (In the special case where $\alpha = 0$ we must, as before, look at higher-order terms in the expansion around the fixed point to determine the nature of the point.)

When there are more than two variables, the eigenvalues must either be real or occur in complex conjugate pairs. Thus again we have eigendirections that simply grow or decay, or that spiral in or out.

We are, however, not done yet. There is a further interesting behavior arising

The imaginary part is also a solution, or any combination of the real and imaginary parts.



The flows around a fixed point whose Jacobian matrix has complex eigenvalues describe a spiral.

in systems with two or more variables that will be important when we come to study networked systems. In addition to fixed points, in some systems one also finds *limit cycles*. A limit cycle is a closed loop in the dynamics such that the system circulates repeatedly through the same set of values. Limit cycles can be treated in many ways rather like fixed points: we can study the dynamics close to a limit cycle by expanding in a small displacement coordinate. Like fixed points, limit cycles can be either attracting or repelling, meaning that points close to them either tend towards the limit cycle or move away from it.

Physically, limit cycles represent stable oscillatory behaviors in systems. We mentioned one such behavior in Section 16.1.6 in our discussion of the SIRS epidemic model. In certain parameter regimes the SIRS model can show “waves” of infection—oscillatory behaviors under which a disease infects a large fraction of the population, who then recover and gain immunity, reducing the number of victims available to the disease and therefore causing the number of cases to drop. When these individuals later lose their immunity, they move back into the susceptible state and become infected again, and another wave starts. Another example of oscillation in a dynamical system is the oscillation of the numbers of predators and prey in a two-species ecosystem represented by the Lotka–Volterra equations [441]. Such oscillations have been famously implicated in the mysterious periodic variation in the populations of hares and lynx recorded by the Hudson Bay Company in Canada during the nineteenth century. We consider limit cycles on networks in Section 17.5.

17.2 DYNAMICS ON NETWORKS

Let us now apply the ideas of the previous section to dynamical systems on networks. First, we need to define the type of systems we are talking about. We will be considering systems in which there are one or more dynamical variables x_i, y_i, \dots on each node i of our network and those variables are coupled together only along the edges of the network. That is, when we write an equation for the time evolution of a variable x_i , the individual terms appearing in that equation each involve only x_i itself, other variables on node i , and/or one or more variables on a node adjacent to i in the network. There are no terms involving variables on non-adjacent nodes and no terms involving variables on more than one adjacent node. (Different terms, however, may involve variables on different adjacent nodes.)

An example of a dynamical system of this type is our equation (16.55) for the probability of infection of a node in the network version of the SI epidemic

model:

$$\frac{dx_i}{dt} = \beta(1 - x_i) \sum_j A_{ij} x_j. \quad (17.33)$$

This equation only has terms involving pairs of variables connected by edges, since these are the only pairs for which A_{ij} is non-zero.

For a network dynamical system with a single variable x_i on each node, the general form of the equation for x_i is

$$\frac{dx_i}{dt} = f_i(x_i) + \sum_j A_{ij} g_{ij}(x_i, x_j), \quad (17.34)$$

where we have separated terms that involve variables on adjacent nodes from those that do not. You can think of $f_i(x_i)$ as specifying the intrinsic dynamics of a node—it tells us how the variable x_i would evolve in the absence of any connections between nodes, i.e., if $A_{ij} = 0$ for all i, j . Conversely, $g_{ij}(x_i, x_j)$ describes the contribution from the connections—it represents the coupling between variables on different nodes.

Note that we have specified different functions f_i and g_{ij} for each node or pair of nodes, so the dynamics obeyed by each node can be different. In many cases, however, when each of the nodes represents the same thing—such as a person in the case of an epidemic model—the dynamics for each node will be the same, or at least similar enough that we can ignore any differences. In such cases, the functions in Eq. (17.34) become the same for all nodes and the equation takes the form

$$\frac{dx_i}{dt} = f(x_i) + \sum_j A_{ij} g(x_i, x_j). \quad (17.35)$$

In the examples in this chapter we will assume that this is the case. We will also assume that the network is undirected so that A_{ij} is symmetric—if x_i is affected by x_j then x_j is similarly affected by x_i . (Note, however, that we do not assume that the function g is symmetric in its arguments; in general $g(u, v) \neq g(v, u)$.) Again, the SI model of Eq. (17.33) is an example of a system of this kind, one in which $f(x) = 0$ and $g(x_i, x_j) = \beta(1 - x_i)x_j$.

17.2.1 LINEAR STABILITY ANALYSIS

Let us apply the tools of linear stability analysis to Eq. (17.35). Suppose we are able to find a fixed point $\{x_i^*\}$ of Eq. (17.35) by solving the simultaneous equations

$$f(x_i^*) + \sum_j A_{ij} g(x_i^*, x_j^*) = 0 \quad (17.36)$$

for all i . Note that finding a fixed point in this case means finding a value $x_i = x_i^*$ for every node i —the fixed point is the complete set $\{x_i^*\}$. Note also that, in general, the position of the fixed point depends both on the particular dynamical process taking place on the network (via the functions f and g) and on the structure of the network (via the adjacency matrix). If either is changed then the position of the fixed point may also change.

Now we can linearize about this fixed point by writing $x_i = x_i^* + \epsilon_i$, performing a multiple Taylor expansion in all variables simultaneously, and dropping terms at second order in small quantities and higher:

$$\begin{aligned}
 \frac{dx_i}{dt} &= \frac{d\epsilon_i}{dt} = f(x_i^* + \epsilon_i) + \sum_j A_{ij}g(x_i^* + \epsilon_i, x_j^* + \epsilon_j) \\
 &= f(x_i^*) + \epsilon_i \left. \frac{df}{dx} \right|_{x=x_i^*} + \sum_j A_{ij}g(x_i^*, x_j^*) \\
 &\quad + \epsilon_i \sum_j A_{ij} \left. \frac{\partial g(u, v)}{\partial u} \right|_{u=x_i^*, v=x_j^*} + \sum_j A_{ij}\epsilon_j \left. \frac{\partial g(u, v)}{\partial v} \right|_{u=x_i^*, v=x_j^*} + \dots \\
 &= \epsilon_i \left. \frac{df}{dx} \right|_{x=x_i^*} + \epsilon_i \sum_j A_{ij} \left. \frac{\partial g(u, v)}{\partial u} \right|_{u=x_i^*, v=x_j^*} + \sum_j A_{ij}\epsilon_j \left. \frac{\partial g(u, v)}{\partial v} \right|_{u=x_i^*, v=x_j^*} + \dots,
 \end{aligned} \tag{17.37}$$

where we have used Eq. (17.36) in the last line.

If we know the position of the fixed point, then the derivatives in these expressions are simply numbers. For convenience, let us write

$$\alpha_i = \left. \frac{df}{dx} \right|_{x=x_i^*}, \tag{17.38a}$$

$$\beta_{ij} = \left. \frac{\partial g(u, v)}{\partial u} \right|_{u=x_i^*, v=x_j^*}, \tag{17.38b}$$

$$\gamma_{ij} = \left. \frac{\partial g(u, v)}{\partial v} \right|_{u=x_i^*, v=x_j^*}. \tag{17.38c}$$

Then Eq. (17.37) becomes

$$\frac{d\epsilon_i}{dt} = \left[\alpha_i + \sum_j \beta_{ij}A_{ij} \right] \epsilon_i + \sum_j \gamma_{ij}A_{ij}\epsilon_j, \tag{17.39}$$

which we can write in matrix form as

$$\frac{d\epsilon}{dt} = \mathbf{M}\epsilon, \tag{17.40}$$

where \mathbf{M} is the matrix with elements

$$M_{ij} = \delta_{ij} \left[\alpha_i + \sum_k \beta_{ik} A_{ik} \right] + \gamma_{ij} A_{ij}, \quad (17.41)$$

and δ_{ij} is the Kronecker delta.

We can solve Eq. (17.40) by writing ϵ as a linear combination of the eigenvectors of \mathbf{M} , specifically the *right* eigenvectors, since \mathbf{M} is in general not symmetric:

$$\epsilon(t) = \sum_r c_r(t) \mathbf{v}_r, \quad (17.42)$$

so that Eq. (17.40) becomes

$$\sum_r \frac{dc_r}{dt} \mathbf{v}_r = \mathbf{M} \sum_r c_r(t) \mathbf{v}_r = \sum_r \mu_r c_r(t) \mathbf{v}_r, \quad (17.43)$$

where μ_r is the eigenvalue corresponding to the eigenvector \mathbf{v}_r . Comparing terms in each eigenvector we then have

$$\frac{dc_r}{dt} = \mu_r c_r(t), \quad (17.44)$$

which implies that

$$c_r(t) = c_r(0) e^{\mu_r t}. \quad (17.45)$$

Immediately we see that if the real parts of all of the eigenvalues μ_r are negative, then $c_r(t)$ —and hence ϵ —is decaying in time for all r and our fixed point will be attracting. If the real parts are all positive the fixed point will be repelling. And if some are positive and some are negative then the fixed point is a saddle, although, as previously, this is perhaps best looked at as a form of repelling fixed point: the flows near a saddle have at least one repelling direction, which means that a system starting in the vicinity of such a point will not in general stay near it, regardless of whether the other directions are attracting or not.

17.2.2 SPECIAL CASES

Let us look at some special cases of the general formalism of Section 17.2.1. A particularly simple case is when the fixed point is *symmetric*, meaning that x_i^* has the same value for every i , so that $x_i^* = x^*$ for some x^* . This occurs in the SI model for instance—there is a fixed point at $x_i^* = 0$ for all i and another at $x_i^* = 1$.

For a symmetric fixed point, the fixed point equation, Eq. (17.36), becomes

$$f(x^*) + \sum_j A_{ij} g(x^*, x^*) = f(x^*) + k_i g(x^*, x^*) = 0, \quad (17.46)$$

where k_i is the degree of node i and we have made use of $k_i = \sum_j A_{ij}$ (see Eq. (6.12)). There are only two ways this equation can be satisfied for all i : either all nodes must have the same degree or $g(x^*, x^*) = 0$. Since the former is not really realistic—few networks of interest have all degrees the same—let us concentrate on the latter and assume that

$$g(x^*, x^*) = 0. \quad (17.47)$$

Again the SI model provides an example of this type of behavior: the coupling function g is of the form $\beta(1-x)x$, which is zero at the two fixed points at $x = 0, 1$.

Equations (17.46) and (17.47) together imply also that $f(x^*) = 0$ and hence the fixed point x^* is the same in this case as the fixed point for the “intrinsic” dynamics of a node: it falls at the same place as it would if there were no connections between nodes at all. The position of the fixed point is also independent of the network structure in this case, a point that will shortly be important.

For a symmetric fixed point, the quantities α_i , β_{ij} , and γ_{ij} defined in Eq. (17.38) become

$$\alpha_i = \alpha = \left. \frac{\partial f}{\partial x} \right|_{x=x^*}, \quad (17.48a)$$

$$\beta_{ij} = \beta = \left. \frac{\partial g(u, v)}{\partial u} \right|_{u, v=x^*}, \quad (17.48b)$$

$$\gamma_{ij} = \gamma = \left. \frac{\partial g(u, v)}{\partial v} \right|_{u, v=x^*}. \quad (17.48c)$$

Then Eq. (17.39) becomes

$$\frac{d\epsilon_i}{dt} = (\alpha + \beta k_i) \epsilon_i + \gamma \sum_j A_{ij} \epsilon_j. \quad (17.49)$$

The situation simplifies further if the coupling function $g(x_i, x_j)$ depends only on x_j and not on x_i , i.e., if x_i obeys an equation of the form $dx_i/dt = f(x_i) + \sum_j A_{ij} g(x_j)$. Then $\beta = 0$ and

$$\frac{d\epsilon_i}{dt} = \alpha \epsilon_i + \gamma \sum_j A_{ij} \epsilon_j, \quad (17.50)$$

which we can write in matrix form as

$$\frac{d\epsilon}{dt} = (\alpha\mathbf{I} + \gamma\mathbf{A})\epsilon. \quad (17.51)$$

As in the general case, the fixed point will be stable if and only if all of the eigenvalues of the matrix $\alpha\mathbf{I} + \gamma\mathbf{A}$ are negative. (The matrix is symmetric, so the eigenvalues are all real.) Let \mathbf{v}_r be the eigenvector of the adjacency matrix with eigenvalue κ_r . Then

$$(\alpha\mathbf{I} + \gamma\mathbf{A})\mathbf{v}_r = \alpha\mathbf{I}\mathbf{v}_r + \gamma\mathbf{A}\mathbf{v}_r = \alpha\mathbf{v}_r + \gamma\kappa_r\mathbf{v}_r = (\alpha + \gamma\kappa_r)\mathbf{v}_r. \quad (17.52)$$

Hence \mathbf{v}_r is also an eigenvector of $\alpha\mathbf{I} + \gamma\mathbf{A}$, but with eigenvalue $\alpha + \gamma\kappa_r$. If all eigenvalues are to be negative, we require that

$$\alpha + \gamma\kappa_r < 0 \quad (17.53)$$

for all r and from this we can deduce a number of things. First of all it implies that $\alpha < -\gamma\kappa_r$ for all r . The adjacency matrix always has both positive and negative eigenvalues (a result that we will prove in Section 17.4), which means that for this inequality to be satisfied for all r we must have $\alpha < 0$ (regardless of whether γ is positive or negative). If $\alpha > 0$ then the fixed point is never stable.

Second, we can rearrange Eq. (17.53) to give

$$\kappa_r < -\alpha/\gamma \quad \text{if } \gamma > 0, \quad (17.54a)$$

$$\kappa_r > -\alpha/\gamma \quad \text{if } \gamma < 0, \quad (17.54b)$$

for all r . Note, however, that if Eq. (17.54a) is satisfied for the most positive eigenvalue κ_1 of the adjacency matrix, then it is necessarily satisfied by all the other eigenvalues as well. Similarly, if Eq. (17.54b) is satisfied for the most negative eigenvalue κ_n then it is satisfied by all others. Thus the system will be stable when

$$\kappa_1 < -\alpha/\gamma \quad \text{if } \gamma > 0, \quad (17.55a)$$

$$\kappa_n > -\alpha/\gamma \quad \text{if } \gamma < 0. \quad (17.55b)$$

Alternatively, we can take reciprocals of these conditions and combine them into a single statement:

$$\frac{1}{\kappa_n} < -\frac{\gamma}{\alpha} < \frac{1}{\kappa_1}. \quad (17.56)$$

For any values of the quantities α and γ , the system will be stable if and only if this expression is satisfied. If we want, we can fill in the explicit values of α and γ from Eq. (17.48) thus:

$$\frac{1}{\kappa_n} < -\left[\frac{dg}{dx} / \frac{df}{dx}\right]_{x=x^*} < \frac{1}{\kappa_1}, \quad (17.57)$$

where we have written g as a function of a single variable since, by hypothesis, it only depends on one argument in this case.

Equation (17.57) is sometimes called a *master stability condition*. It has a special form: note that κ_1 and κ_n depend only on the structure of the network and not on anything about the dynamics, while the derivatives of f and g depend only the nature of the dynamics and not on the network structure. Thus Eq. (17.57) effectively gives us a single condition that must be satisfied by any type of dynamics if that dynamics is to be stable on a given network. Or conversely, it gives a condition on the network structure, via the largest and smallest eigenvalues, that guarantees stability of the fixed point for a given type of dynamics.

Another case where we can derive a master stability condition is the case in which the coupling function g depends on its two arguments according to $g(x_i, x_j) = g(x_i) - g(x_j)$. A physicist might think of this as a “spring-like” interaction—if $g(x)$ were a simple linear function of its argument, then x_i and x_j would act upon one another like two masses coupled by a spring, exerting forces that depend on the difference of their positions. More generally, $g(x)$ is non-linear and we have a non-linear spring.

For this choice of coupling, and still assuming a symmetric fixed point, we have $g(x^*, x^*) = 0$ as before and hence also $f(x^*) = 0$ (via Eq. (17.46)), and the quantities defined in Eq. (17.38) become

$$\alpha_i = \alpha = \left. \frac{df}{dx} \right|_{x=x^*}, \quad (17.58a)$$

$$\beta_{ij} = \beta = \left. \frac{dg}{dx} \right|_{x=x^*}, \quad (17.58b)$$

$$\gamma_{ij} = -\beta. \quad (17.58c)$$

Then Eq. (17.39) becomes

$$\frac{d\epsilon_i}{dt} = (\alpha + \beta k_i)\epsilon_i - \beta \sum_j A_{ij}\epsilon_j = \alpha\epsilon_i + \beta \sum_j (k_i\delta_{ij} - A_{ij})\epsilon_j, \quad (17.59)$$

or in matrix form

$$\frac{d\epsilon}{dt} = (\alpha\mathbf{I} + \beta\mathbf{L})\epsilon, \quad (17.60)$$

where \mathbf{L} is the graph Laplacian matrix, the matrix with elements

$$L_{ij} = k_i\delta_{ij} - A_{ij}. \quad (17.61)$$

We encountered the graph Laplacian previously in Section 6.14 (see Eq. (6.28) on page 143) and in Section 14.7.4.

Equation (17.60) is of the same form as Eq. (17.51) except that the adjacency matrix is replaced by the Laplacian. Thus, following the same argument that led to Eq. (17.53), we can see that the fixed point will be stable if and only if the eigenvalues λ_r of the Laplacian satisfy

$$\alpha + \beta\lambda_r < 0 \tag{17.62}$$

for all r .

As shown in Section 6.14.5, the smallest eigenvalue of the Laplacian matrix is always zero, and hence Eq. (17.62), when applied to the smallest eigenvalue, implies again that $\alpha < 0$, or equivalently

$$\left. \frac{df}{dx} \right|_{x=x^*} < 0, \tag{17.63}$$

if the system is to be stable. Assuming this condition is satisfied, then since all eigenvalues of the Laplacian are non-negative it follows that $1/\lambda_r > -\beta/\alpha$ for stability, regardless of the sign of β . Furthermore, if this condition is true for the largest eigenvalue, conventionally denoted λ_n , then it is true for all smaller eigenvalues as well, so the requirement for stability can be reduced to the single condition $1/\lambda_n > -\beta/\alpha$, or

$$\frac{1}{\lambda_n} > - \left[\frac{dg}{dx} / \frac{df}{dx} \right]_{x=x^*}, \tag{17.64}$$

along with the constraint in Eq. (17.63).

Equation (17.64) is another example of a master stability condition, and again it neatly separates questions of dynamics from questions of network structure. The structure appears only on the left of the inequality, via the largest eigenvalue of the graph Laplacian, and the dynamics appears only on the right, via derivatives of the functions f and g .

Apart from establishing the requirements for the stability of a fixed point, master stability conditions are also of interest in the study of bifurcations—situations in which a fixed point loses stability as the parameters of a system change. If we vary the functions f and g , for example, then we can cause a fixed point that initially satisfies a condition like (17.64) to stop satisfying it and so become unstable. In practice, this means that the system will suddenly change its behavior as it passes through the point where $1/\lambda_n = -\beta/\alpha$. At one moment it will be sitting happily at its stable fixed point, going nowhere, and at the next, as that point becomes unstable, it will start moving, gathering speed exponentially, and quite likely wind up in some completely different state far from where it started, as it falls into the basin of attraction of a different stable fixed point or limit cycle. We will see some examples of behavior of this kind shortly.

17.2.3 AN EXAMPLE

As an example, consider the following simple model of “gossip,” or diffusion of an idea or fad across a social network. Suppose some new idea is circulating through a community and x_i represents the amount person i is talking about it. The value of x_i is governed by an equation of the form (17.35), which we repeat here for convenience:

$$\frac{dx_i}{dt} = f(x_i) + \sum_j A_{ij}g(x_i, x_j). \quad (17.65)$$

We will put

$$f(x) = a(1 - x) \quad (17.66)$$

with $a > 0$, which means that the intrinsic dynamics of a single node has a stable fixed point at $x^* = 1$ —each person has an intrinsic tendency to talk this much about the latest craze, regardless of whether their friends want to hear about it. For the interaction term we will assume that people tend either to copy their friends or to differ from them: if their friends are talking about whatever it is more than they are, then over time they themselves will talk about it either more or less, depending on the parameters of the model. We represent this by putting $g(x_i, x_j) = g(x_i) - g(x_j)$, as in Section 17.2.2, with

$$g(x) = \frac{bx}{1 + x}. \quad (17.67)$$

When $b < 0$ this implies that $g(x_i, x_j) > 0$ when $x_j > x_i$, so that people tend to copy their friends; when $b > 0$ we have $g(x_i, x_j) < 0$ and they try to differentiate themselves. However, the function saturates when $x \gg 1$, so that beyond some point it makes no difference if your friends shout louder.

Now we can apply the general formalism developed in previous sections. The symmetric fixed point for the model is at $x_i = 1$ for all i . At this point everyone is talking about the topic du jour with equal enthusiasm. This fixed point, however, is stable only provided Eqs. (17.63) and (17.64) are satisfied. Equation (17.63) is always satisfied, given that $a > 0$. Equation (17.64) implies that for stability we must have

$$\frac{1}{\lambda_n} > \frac{b}{4a}. \quad (17.68)$$

This condition is always satisfied if $b < 0$ since λ_n cannot be negative, so the fixed point is always stable when people are copying their friends. If people are contrary, however, and try to be different from their friends, then $b > 0$ and the fixed point becomes unstable when $b > 4a/\lambda_n$.

And what happens when the fixed point becomes unstable? There are no other symmetric fixed points for this particular system, since there are no other values that give $f(x) = 0$ (which is a requirement for our symmetric fixed point). So the system cannot switch to another symmetric fixed point. One possibility is that the variables might diverge to $\pm\infty$, and this happens in some systems, but not in this one, where the form of $f(x)$ prevents it. Another possibility is that the system might begin to oscillate, or even enter a chaotic regime in which it meanders around in pseudorandom fashion indefinitely. In the present case, however, it does something simpler. It moves to a *non-symmetric* fixed point, one in which the fixed-point values of the variables x_i are not all equal. In other words if the influence between neighboring individuals in the network is strong and contrary, if people really want to be different from their neighbors, then they will eventually start doing things their own way, but only for sufficiently strong interactions—people will put up with a weak tendency towards nonconformity without developing idiosyncrasies.

17.3 DYNAMICS WITH MORE THAN ONE VARIABLE PER NODE

Our developments so far have assumed that there is only a single variable x_i on each node i of the network. Many systems, however, have more than one variable per node. The epidemiological examples of Chapter 16, for instance, all have two or more— s, x, r , and so forth.

Consider a system with an arbitrary number of variables x_1^i, x_2^i, \dots on each node i , but let us assume that we have the same number of variables on every node and that, as before, every node obeys equations of the same form. For convenience let us write the set of variables on a single node as a vector $\mathbf{x}^i = (x_1^i, x_2^i, \dots)$ and then write the equations governing their time evolution as

$$\frac{d\mathbf{x}^i}{dt} = \mathbf{f}(\mathbf{x}^i) + \sum_j A_{ij} \mathbf{g}(\mathbf{x}^i, \mathbf{x}^j). \quad (17.69)$$

Note that the functions f and g , representing the intrinsic dynamics and the coupling, have now become vector functions \mathbf{f} and \mathbf{g} of vector arguments, with the same number of elements as \mathbf{x} .

Following the same line of reasoning as before, we can study the stability of a fixed point $\mathbf{x}^i = \mathbf{x}^*$ by writing $\mathbf{x}^i = \mathbf{x}^* + \boldsymbol{\epsilon}^i$ and performing a Taylor expansion. The resulting linearized equation for the evolution of the μ th component of $\boldsymbol{\epsilon}^i$ is then

$$\begin{aligned}
 \frac{d\epsilon_\mu^i}{dt} &= \left[\epsilon_1^i \frac{\partial f_\mu(\mathbf{x})}{\partial x_1} + \epsilon_2^i \frac{\partial f_\mu(\mathbf{x})}{\partial x_2} + \dots \right]_{\mathbf{x}=\mathbf{x}^*} \\
 &\quad + \sum_j A_{ij} \left[\epsilon_1^i \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial u_1} + \epsilon_2^i \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial u_2} + \dots + \epsilon_1^j \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial v_1} + \epsilon_2^j \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial v_2} + \dots \right]_{\mathbf{u}, \mathbf{v}=\mathbf{x}^*} \\
 &= \sum_\nu \left[\epsilon_\nu^i \frac{\partial f_\mu(\mathbf{x})}{\partial x_\nu} \Big|_{\mathbf{x}=\mathbf{x}^*} + k_i \epsilon_\nu^i \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial u_\nu} \Big|_{\mathbf{u}, \mathbf{v}=\mathbf{x}^*} + \sum_j A_{ij} \epsilon_\nu^j \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial v_\nu} \Big|_{\mathbf{u}, \mathbf{v}=\mathbf{x}^*} \right], \tag{17.70}
 \end{aligned}$$

where f_μ and g_μ represent the μ th components of \mathbf{f} and \mathbf{g} .

As before, the derivatives in this expression are simply constants, and for convenience let us define

$$\alpha_{\mu\nu} = \frac{\partial f_\mu(\mathbf{x})}{\partial x_\nu} \Big|_{\mathbf{x}=\mathbf{x}^*}, \tag{17.71a}$$

$$\beta_{\mu\nu} = \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial u_\nu} \Big|_{\mathbf{u}, \mathbf{v}=\mathbf{x}^*}, \tag{17.71b}$$

$$\gamma_{\mu\nu} = \frac{\partial g_\mu(\mathbf{u}, \mathbf{v})}{\partial v_\nu} \Big|_{\mathbf{u}, \mathbf{v}=\mathbf{x}^*}, \tag{17.71c}$$

so that

$$\begin{aligned}
 \frac{d\epsilon_\mu^i}{dt} &= \sum_\nu [(\alpha_{\mu\nu} + k_i \beta_{\mu\nu}) \epsilon_\nu^i + \sum_j A_{ij} \gamma_{\mu\nu} \epsilon_\nu^j] \\
 &= \sum_{j\nu} [\delta_{ij} (\alpha_{\mu\nu} + k_i \beta_{\mu\nu}) + A_{ij} \gamma_{\mu\nu}] \epsilon_\nu^j \tag{17.72}
 \end{aligned}$$

where δ_{ij} is the Kronecker delta again.

We can write this equation in the matrix form

$$\frac{d\boldsymbol{\epsilon}}{dt} = \mathbf{M}\boldsymbol{\epsilon}, \tag{17.73}$$

where \mathbf{M} is a matrix whose rows (and columns) are labeled by a double pair of indices (i, μ) and whose elements are

$$M_{i\mu, j\nu} = \delta_{ij} (\alpha_{\mu\nu} + k_i \beta_{\mu\nu}) + A_{ij} \gamma_{\mu\nu}. \tag{17.74}$$

In principle, we can now determine whether the fixed point is stable by examining the eigenvalues of this new matrix. If the real parts of the eigenvalues are all negative, then the fixed point is stable; otherwise it is not. In practice this can be a difficult thing to do in general but, as before, there are some special cases where the calculation simplifies, yielding a master stability condition.

17.3.1 SPECIAL CASES

As before, we consider the case where $\mathbf{g}(\mathbf{x}^i, \mathbf{x}^j)$ depends only on its second argument and not on its first. In this case $\beta_{\mu\nu} = 0$ for all μ, ν and Eq. (17.72) becomes

$$\frac{d\epsilon_{\mu}^i}{dt} = \sum_{j\nu} [\delta_{ij}\alpha_{\mu\nu} + A_{ij}\gamma_{\mu\nu}] \epsilon_{\nu}^j. \quad (17.75)$$

Now let v_r^i be the i th component of the eigenvector \mathbf{v}_r of the adjacency matrix corresponding to eigenvalue κ_r . Let us write

$$\epsilon_{\mu}^i(t) = \sum_r c_{\mu}^r(t) v_r^i. \quad (17.76)$$

This equation expresses the vector of elements ϵ_{μ}^i as a linear combination of eigenvectors in the usual way, but with a separate set of coefficients c_{μ}^r for each dynamical variable μ . Substituting into Eq. (17.75), we get

$$\begin{aligned} \sum_r \frac{dc_{\mu}^r}{dt} v_r^i &= \sum_r \sum_{j\nu} [\delta_{ij}\alpha_{\mu\nu} + A_{ij}\gamma_{\mu\nu}] c_{\nu}^r(t) v_r^j \\ &= \sum_{r\nu} [\alpha_{\mu\nu} + \kappa_r \gamma_{\mu\nu}] c_{\nu}^r(t) v_r^i. \end{aligned} \quad (17.77)$$

Equating terms in the individual eigenvectors on both sides of the equation, we conclude that

$$\frac{dc_{\mu}^r}{dt} = \sum_{\nu} [\alpha_{\mu\nu} + \kappa_r \gamma_{\mu\nu}] c_{\nu}^r(t). \quad (17.78)$$

We can think of this as an equation for a vector $\mathbf{c}^r = (c_1^r, c_2^r, \dots)$ thus:

$$\frac{d\mathbf{c}^r}{dt} = [\boldsymbol{\alpha} + \kappa_r \boldsymbol{\gamma}] \mathbf{c}^r(t), \quad (17.79)$$

where $\boldsymbol{\alpha}$ and $\boldsymbol{\gamma}$ are matrices with elements $\alpha_{\mu\nu}$ and $\gamma_{\mu\nu}$, respectively. This equation expresses the dynamics of the system close to the fixed point as a decoupled set of n separate systems, one for each eigenvalue κ_r of the adjacency matrix. If the fixed point of the system as a whole is to be stable, then each of these individual systems also needs to be stable, meaning that *their* eigenvalues need to be negative, or, more simply, the largest (i.e., most positive) eigenvalue of $\boldsymbol{\alpha} + \kappa_r \boldsymbol{\gamma}$ needs to be negative for every r .

Let us define the function $\sigma(\kappa)$ to be equal to the most positive eigenvalue of the matrix $\boldsymbol{\alpha} + \kappa \boldsymbol{\gamma}$, or the most positive real part in the case where the eigenvalues are complex. Typically, this is an easy function to evaluate numerically. Note

that $\alpha + \kappa\gamma$ has only as many rows and columns as there are variables on each node of the network. If we have three variables on each node, for instance, the matrix has size 3×3 , which is easily diagonalized.

The function $\sigma(\kappa)$ is called a *master stability function*. If our system is to be stable, the master stability function evaluated at the eigenvalue κ_r should be negative for all r :

$$\sigma(\kappa_r) < 0. \quad (17.80)$$

One possible form for the master stability function is shown in Fig. 17.3—it becomes large and positive for κ sufficiently small or sufficiently big, but is negative in some intermediate range $\kappa_{\min} < \kappa < \kappa_{\max}$. In that case, the system is stable provided all eigenvalues κ_r of the adjacency matrix fall in this range. Again this gives us a master stability condition that separates network structure from dynamics. The eigenvalues κ_r are properties solely of the structure, being derived from the adjacency matrix alone, while the limits κ_{\min} and κ_{\max} are properties solely of the dynamics, being derived from the matrices α and γ , which are determined by the derivatives of the functions f and g .

We can perform a similar analysis for other special cases as well. For instance, there is a natural generalization of Eqs. (17.58) to (17.60) to the case of many variables per node. If the interaction between nodes takes the form $\mathbf{g}(\mathbf{x}^i, \mathbf{x}^j) = \mathbf{g}(\mathbf{x}^i) - \mathbf{g}(\mathbf{x}^j)$, then $\gamma_{\mu\nu} = -\beta_{\mu\nu}$ in Eq. (17.71) and

$$\frac{d\epsilon_\mu^i}{dt} = \sum_{j\nu} [\delta_{ij}\alpha_{\mu\nu} + L_{ij}\beta_{\mu\nu}] \epsilon_\nu^j, \quad (17.81)$$

where $L_{ij} = k_i\delta_{ij} - A_{ij}$ is an element of the graph Laplacian. Then the equivalent of Eq. (17.79) is

$$\frac{d\mathbf{c}^r}{dt} = [\alpha + \lambda_r\beta]\mathbf{c}^r(t), \quad (17.82)$$

where λ_r is an eigenvalue of the Laplacian and β is the matrix with elements $\beta_{\mu\nu}$. Again we can define a master stability function $\sigma(\lambda)$ equal to the most positive eigenvalue of $\alpha + \lambda\beta$ (or the most positive real part for complex eigenvalues) and for overall stability of the system this function must be negative at λ_r for all r :

$$\sigma(\lambda_r) < 0. \quad (17.83)$$

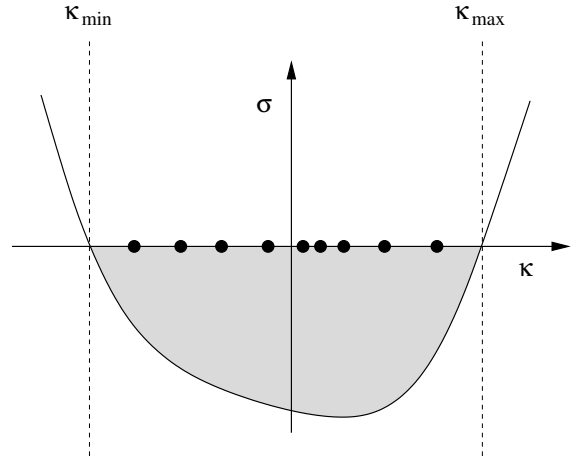


Figure 17.3: A sketch of a master stability function. One possible form for the master stability function $\sigma(\kappa)$ might be as shown here (solid curve), with positive values for large and small κ but negative values in an intermediate range between κ_{\min} and κ_{\max} . If all the eigenvalues of the adjacency matrix (represented by the dots) fall in this intermediate range, then the system is stable.

For suitable forms of the master stability function, this allows us once again to develop a stability criterion that separates structure from dynamics.

17.4 SPECTRA OF NETWORKS

The formalism of the previous sections turns questions about the stability of dynamical systems on networks into questions about the eigenvalue spectra of matrices such as the adjacency matrix or the graph Laplacian. For systems with only a single dynamical variable on each node, master stability conditions such as Eqs. (17.57) and (17.64) tell us whether the system is stable in terms of the largest (most positive) or smallest (most negative) eigenvalues of the appropriate matrix. For systems with more than one variable per node we calculate a master stability function and then the stability of the system depends on whether the function is negative when evaluated at each of the eigenvalues. When the master stability function takes a relatively simple form like that sketched in Fig. 17.3, so that stability requires only that the eigenvalues fall in some specified range, then it is enough to know only the smallest and largest eigenvalues of the matrix to ensure stability—if the smallest and largest fall in the required range, then necessarily all the others do too.

A number of results are known about the spectra of networks, and in particular about the smallest and largest eigenvalues, which allow us to make quite general theoretical statements about stability. For the adjacency matrix, for example, we can derive bounds on the largest eigenvalue as follows.

Consider an undirected network and let \mathbf{x} be an arbitrary real vector of n elements, which we will write as a linear combination of the eigenvectors \mathbf{v}_r of the adjacency matrix \mathbf{A} thus:

$$\mathbf{x} = \sum_r c_r \mathbf{v}_r. \quad (17.84)$$

Then

$$\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\sum_s c_s \mathbf{v}_s^T \mathbf{A} \sum_r c_r \mathbf{v}_r}{\sum_s c_s \mathbf{v}_s^T \sum_r c_r \mathbf{v}_r} = \frac{\sum_{rs} c_s c_r \kappa_r \mathbf{v}_s^T \mathbf{v}_r}{\sum_{rs} c_s c_r \mathbf{v}_s^T \mathbf{v}_r} = \frac{\sum_r c_r^2 \kappa_r}{\sum_r c_r^2} \leq \frac{\sum_r c_r^2 \kappa_1}{\sum_r c_r^2} = \kappa_1, \quad (17.85)$$

where, as before, κ_1 is the largest eigenvalue and we have made use of the fact that $\mathbf{v}_s^T \mathbf{v}_r = \delta_{rs}$.

This inequality is true for any choice of \mathbf{x} . Thus, for instance, if $\mathbf{x} = \mathbf{1} = (1, 1, 1, \dots)$, then

$$\kappa_1 \geq \frac{\mathbf{1}^T \mathbf{A} \mathbf{1}}{\mathbf{1}^T \mathbf{1}} = \frac{2m}{n} = \langle k \rangle. \quad (17.86)$$

So the largest eigenvalue of the adjacency matrix is never less than the average degree of the network. This is a well-known and widely used result, but with a little work we can do better.

If \mathbf{v}_r is the eigenvector corresponding to eigenvalue κ_r then $\mathbf{A}^2\mathbf{v}_r = \kappa_r^2\mathbf{v}_r$. In other words the eigenvalues of \mathbf{A}^2 are the squares of the eigenvalues of \mathbf{A} . Moreover, the eigenvalue of the adjacency matrix with largest magnitude is always positive, meaning it is the most positive eigenvalue κ_1 , not the most negative eigenvalue κ_n . (This is a consequence of the Perron–Frobenius theorem—see footnote 2 on page 161.) Thus the largest eigenvalue of \mathbf{A}^2 is κ_1^2 and the equivalent of Eq. (17.85) for the matrix \mathbf{A}^2 is

$$\frac{\mathbf{x}^T \mathbf{A}^2 \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \leq \kappa_1^2. \quad (17.87)$$

Again setting $\mathbf{x} = \mathbf{1} = (1, 1, 1, \dots)$, we then find that

$$\kappa_1^2 \geq \frac{\sum_i k_i^2}{n} = \langle k^2 \rangle, \quad (17.88)$$

and hence

$$\kappa_1 \geq \sqrt{\langle k^2 \rangle}, \quad (17.89)$$

which gives us a bound on κ_1 different from that of (17.86)

The variance $\sigma^2 = \langle k^2 \rangle - \langle k \rangle^2$ of the degree distribution cannot be negative, so $\langle k^2 \rangle \geq \langle k \rangle^2$ always, or equivalently $\sqrt{\langle k^2 \rangle} \geq \langle k \rangle$. Hence (17.89) is always a better bound (or at least no worse) than (17.86). In particular, in the common case of a network with a right-skewed degree distribution (see Section 10.4), the variance σ^2 can become very large so that $\langle k^2 \rangle \gg \langle k \rangle^2$ and (17.89) imposes a much more stringent bound on the leading eigenvalue than does (17.86). Indeed, in networks with a power-law degree distribution with exponent $\alpha \leq 3$, the value of $\langle k^2 \rangle$ formally diverges as we have seen (Section 10.4.2), so we would expect κ_1 to diverge too.

A different bound on the leading eigenvalue can be derived as follows. Suppose that node v is the node of highest degree in the network, with degree k_{\max} , and let us apply Eq. (17.85) again with the elements of \mathbf{x} chosen thus:

$$x_i = \begin{cases} \sqrt{k_{\max}} & \text{if } i = v, \\ 1 & \text{if } A_{iv} = 1, \\ 0 & \text{otherwise.} \end{cases} \quad (17.90)$$

Then

$$\sum_j A_{ij} x_j \geq \begin{cases} k_{\max} & \text{if } i = v \\ \sqrt{k_{\max}} & \text{if } A_{iv} = 1 \\ 0 & \text{otherwise} \end{cases} = \sqrt{k_{\max}} x_i. \quad (17.91)$$

(This result is non-trivial and you may find it helpful to work through each of the three cases to convince yourself that it is indeed correct.)

Multiplying both sides of Eq. (17.91) by x_i and summing over i we now get $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq \sqrt{k_{\max}} \mathbf{x}^T \mathbf{x}$ or, using Eq. (17.85),

$$\kappa_1 \geq \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \geq \sqrt{k_{\max}}. \quad (17.92)$$

Thus, the largest eigenvalue of the adjacency matrix is never less than the square root of the largest degree.

Both (17.89) and (17.92) are useful bounds: which one is better depends on the specifics of the degree distribution. Both of them, however, imply that by increasing the degrees of some or all of the nodes in a network, we can increase the maximum eigenvalue also. In a system with a master stability function like that depicted in Fig. 17.3, this will eventually cause the system to become unstable.

Unfortunately, there is no obvious equivalent of these results for the smallest (most negative) eigenvalue κ_n of the adjacency matrix. Equation (17.85) does have a straightforward generalization:

$$\frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\sum_r c_r^2 \kappa_r}{\sum_r c_r^2} \geq \frac{\sum_r c_r^2 \kappa_n}{\sum_r c_r^2} = \kappa_n \quad (17.93)$$

for any real vector \mathbf{x} . But there are no longer simple choices for \mathbf{x} that give useful general limits on the eigenvalue.

We can, however, see that for a network with no self-edges, so that the diagonal elements of the adjacency matrix are all zero, the trace of the matrix is also zero and hence so is the sum of the eigenvalues. It then follows that the matrix must have both positive and negative eigenvalues, unless the network has no edges in it at all, in which case all eigenvalues are zero. (We used this result previously in Section 17.2.2.) Moreover, the sum of the absolute values of the negative eigenvalues must equal the sum of the positive ones. Without a knowledge of how many eigenvalues there are of either sign, however, one cannot turn this observation into a useful limit on κ_n .

Other results for the eigenvalues of the adjacency matrix can be derived for specific models of networks. For example, Chung *et al.* [104] have shown for the configuration model that the expected value of the largest eigenvalue in the limit of large network size is

$$\kappa_1 = \frac{\langle k^2 \rangle}{\langle k \rangle}. \quad (17.94)$$

In many cases this gives values of κ_1 considerably above the limits set by Eqs. (17.89) and (17.92).

One can also derive results for eigenvalues of the Laplacian. The smallest eigenvalue of the Laplacian is simple—it is always zero. For large networks the largest eigenvalue λ_n can be shown to lie in the range [22]

$$k_{\max} \leq \lambda_n \leq 2k_{\max}, \quad (17.95)$$

where k_{\max} is again the highest degree in the network. The range spanned by (17.95) appears to be a relatively large one but in fact this result tells us a lot, ensuring again that the largest eigenvalue is guaranteed to increase if the highest degree in the network increases sufficiently.

17.5 SYNCHRONIZATION

A topic closely related to dynamical stability on networks is synchronization. Many systems of scientific interest can be thought of as oscillators of one sort or another and a range of interesting phenomena arise when multiple oscillators are linked together in a network. The beating of a human heart, the synchronized clapping of a large audience, the ticking of clocks, flashing of fireflies, or the pathologically synchronized firing of brain cells during an epileptic attack can all be modeled as networks of oscillators coupled in such a way that the oscillators synchronize, meaning that they oscillate in time with one another.

Synchronized oscillators correspond, in dynamical systems terms, to a limit cycle of the overall dynamics, a periodic motion of the variables of the system (see Section 17.1.1). Like fixed points, limit cycles can be stable or unstable, attracting or repelling, depending on whether small perturbations away from the periodic behavior tend to grow or decay over time. Only stable limit cycles give rise to real-world synchronization: unstable states by definition do not last long and are soon replaced by other behaviors.

The mathematics of whether synchronized states are stable is similar to that for fixed points. We can treat oscillating systems using the same kind of equations that we used previously in our study of fixed points, such as Eq. (17.69), which we repeat here for convenience:

$$\frac{d\mathbf{x}^i}{dt} = \mathbf{f}(\mathbf{x}^i) + \sum_j A_{ij} \mathbf{g}(\mathbf{x}^i, \mathbf{x}^j). \quad (17.96)$$

A crucial condition is that we now must have at least two variables on each node of the network, meaning that \mathbf{x}_i must have two or more elements. It is not possible to get oscillating behavior on individual nodes from just a single variable. The first-order differential equation $dx/dt = f(x)$ has no oscillating

solutions, no matter what the form of $f(x)$. On the other hand, the pair of first-order equations

$$\frac{dx}{dt} = y, \quad \frac{dy}{dt} = -x, \quad (17.97)$$

does have oscillating solutions, as we can see by differentiating the first equation and substituting from the second to get

$$\frac{d^2x}{dt^2} = -x, \quad (17.98)$$

which is a standard simple harmonic oscillator.²

So suppose we have a dynamical system of the form (17.96) with at least two variables on each node. First, consider the case where there are no interactions between nodes, meaning that there are no edges in the network and $A_{ij} = 0$ for all i, j . Then Eq. (17.96) becomes simply

$$\frac{dx^i}{dt} = \mathbf{f}(x^i). \quad (17.99)$$

Let us assume that this equation has an oscillating solution $\mathbf{s}(t)$, a limit cycle, meaning that

$$\frac{d\mathbf{s}}{dt} = \mathbf{f}(\mathbf{s}), \quad (17.100)$$

where $\mathbf{s}(t)$ is periodic with some period τ so that $\mathbf{s}(t + \tau) = \mathbf{s}(t)$ for all t .

Given that $\mathbf{s}(t)$ is a solution of (17.99), so also is $\mathbf{s}(t + \phi)$ for any value of ϕ . This means that different nodes need not be executing the motion at the same time: we can have solutions of the form

$$\mathbf{x}^i = \mathbf{s}(t + \phi_i), \quad (17.101)$$

with (potentially) a different value of ϕ_i at every node i . In other words all nodes are following basically the same dynamics, but they do so on their own time, each at a different point around the cycle defined by $\mathbf{s}(t)$. The fireflies are

²An alternative approach, which allows us to represent an oscillator with a single first-order equation, is to model not the oscillator itself but its phase $\theta(t)$. That is, we define a variable $x = \sin \theta$ or similar, and then write a differential equation for θ , such as $d\theta/dt = \omega$, where ω is a constant. This clearly has solutions $\theta = \omega t + \phi$ where ϕ is an integration constant, and hence $x = \sin(\omega t + \phi)$, which oscillates with angular frequency ω . Coupling between oscillators on different nodes i, j of a network is then written as a coupling between their phases θ_i and θ_j . This approach places some limitations on the dynamics—not all oscillating functions can be written as $\sin \theta$ for some θ and not all interactions can be written in terms of phases—but it can make the equations simpler. The classic example of a model of this kind is the Kuramoto synchronization model—see Strogatz [441] for a discussion.

all flashing at the same rate, but not at the same time. The audience members are all clapping at the same speed but they are not synchronized with each other.

Now let us reintroduce the network and the interactions between nodes that it represents. These interactions, it turns out, can be enough to make the nodes synchronize with one another. The simplest case, and the one we will focus on here, is when the interaction function \mathbf{g} takes the form $\mathbf{g}(\mathbf{x}, \mathbf{y}) = \mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})$. Then there exists a solution of the equations such that $\mathbf{x}^i = \mathbf{s}(t)$ for all i and all t , as we can show by substituting this form into Eq. (17.96), which recovers Eq. (17.100). Thus if we set all $\mathbf{x}^i = \mathbf{s}(t)$ at any time t they will remain equal forever afterward. This is what we mean by the synchronized state. All nodes are following the same dynamics and they are now perfectly in time with each other: each is at the same point around the limit cycle at every moment. All the fireflies are flashing in time; all the audience members are clapping together. On the other hand, it is no longer true in general that states of the form (17.101) are solutions to Eq. (17.96). Thus, by introducing interactions between the nodes we have removed the unsynchronized solutions but not the synchronized ones.

Now we can study the stability of the synchronized state in a manner similar to the way we studied the stability of fixed points in previous sections. If we make a small perturbation around the synchronized state—if one audience member is clapping out of time, say—will that perturbation die away over time or will it grow? Following an argument analogous to that of Section 17.3, we write $\mathbf{x}^i(t) = \mathbf{s}(t) + \boldsymbol{\epsilon}^i(t)$, where $\boldsymbol{\epsilon}^i(t)$ is a small quantity. Substituting into Eq. (17.96) and performing a Taylor expansion, we arrive at the equivalent of Eq. (17.81):

$$\frac{d\boldsymbol{\epsilon}_\mu^i}{dt} = \sum_{jv} [\delta_{ij}\alpha_{\mu v}(t) + L_{ij}\beta_{\mu v}(t)] \boldsymbol{\epsilon}_v^j, \quad (17.102)$$

where $L_{ij} = k_i\delta_{ij} - A_{ij}$ is an element of the graph Laplacian and

$$\alpha_{\mu v}(t) = \left. \frac{\partial f_\mu}{\partial x_v} \right|_{\mathbf{x}=\mathbf{s}(t)}, \quad \beta_{\mu v}(t) = \left. \frac{\partial g_\mu}{\partial x_v} \right|_{\mathbf{x}=\mathbf{s}(t)}, \quad (17.103)$$

with f_μ and g_μ representing the μ th components of \mathbf{f} and \mathbf{g} as before. Note that $\alpha_{\mu v}$ and $\beta_{\mu v}$ are now time-dependent, since the derivatives are evaluated at $\mathbf{s}(t)$, which is itself time-dependent. Moreover, they are also periodic with the same period as $\mathbf{s}(t)$.

Now consider the n values $\boldsymbol{\epsilon}_\mu^i$ for one particular value of μ and all i , and let us put them together in the form of an n -element vector $\boldsymbol{\epsilon}_\mu$. In terms of this

vector, Eq. (17.102) can be written

$$\frac{d\epsilon_\mu}{dt} = \sum_\nu [\alpha_{\mu\nu}(t)\mathbf{I} + \beta_{\mu\nu}(t)\mathbf{L}]\epsilon_\nu. \quad (17.104)$$

Let us write ϵ_μ as a linear combination of the eigenvectors \mathbf{v}_r of the Laplacian:

$$\epsilon_\mu(t) = \sum_r c_\mu^r(t)\mathbf{v}_r, \quad (17.105)$$

in which case Eq. (17.104) becomes

$$\begin{aligned} \sum_r \frac{dc_\mu^r}{dt}\mathbf{v}_r &= \sum_\nu [\alpha_{\mu\nu}(t)\mathbf{I} + \beta_{\mu\nu}(t)\mathbf{L}] \sum_r c_r(t)\mathbf{v}_r \\ &= \sum_{\nu r} [\alpha_{\mu\nu}(t) + \lambda_r\beta_{\mu\nu}(t)]c_\nu^r(t)\mathbf{v}_r, \end{aligned} \quad (17.106)$$

where λ_r is the r th eigenvalue of the Laplacian, as previously. Comparing terms in \mathbf{v}_r we then have

$$\frac{dc_\mu^r}{dt} = \sum_\nu [\alpha_{\mu\nu}(t) + \lambda_r\beta_{\mu\nu}(t)]c_\nu^r(t). \quad (17.107)$$

Or we can think of this as a matrix equation for the vector $\mathbf{c}^r = (c_1^r, c_2^r, \dots)$, of the form

$$\frac{d\mathbf{c}^r}{dt} = [\boldsymbol{\alpha}(t) + \lambda_r\boldsymbol{\beta}(t)]\mathbf{c}^r(t). \quad (17.108)$$

This differs from Eq. (17.82) for fixed points in that the matrices $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are now time-dependent. This, however, is enough to make it much harder to find a solution. We can no longer write a simple exponential solution as we did in Eq. (17.45). A complete solution is still possible, but the calculations are complicated and can usually only be done numerically [441]. Here we take a simpler approach: rather than finding a complete solution we will merely derive a sufficient condition for stability.

Suppose we observe the state of the system at some time t and at that time we freeze $\boldsymbol{\alpha}(t)$ and $\boldsymbol{\beta}(t)$ at their current values. Then, since these quantities are now constant, we can apply the same arguments we used for fixed points in Section 17.2.2 and define a master stability function $\sigma_t(\lambda)$ equal to the most positive eigenvalue of $\boldsymbol{\alpha}(t) + \lambda\boldsymbol{\beta}(t)$, or the most positive real part in the case of complex eigenvalues. Then if $\sigma(\lambda_r) < 0$ for all r the system is stable—the perturbation ϵ_μ^i will die away.

Now in fact the values of $\boldsymbol{\alpha}(t)$ and $\boldsymbol{\beta}(t)$ are not constant but change over time as we have seen, so this calculation only applies instantaneously at one

particular time. When we freeze $\alpha(t)$ and $\beta(t)$, the corresponding stability calculation tells us whether the perturbations are dying away at that moment; at a different moment they could grow again. However, if they are dying away at every moment and never growing then necessarily the system will be stable. In other words, if $\sigma_t(\lambda_r) < 0$ for all r and all t then the system is stable. Note, however, that since $\alpha(t)$ and $\beta(t)$ are periodic with the same period τ as $\mathbf{s}(t)$, then so is $\sigma_t(\lambda)$, and hence we only need to check whether $\sigma_t(\lambda_r) < 0$ for values of t spanning one period—from $t = 0$ to $t = \tau$ say. If this condition is satisfied then the system is stable for all time.

So let us define a new master stability function

$$\sigma(\lambda) = \max_{t \in [0, \tau]} \sigma_t(\lambda). \quad (17.109)$$

That is, the function is equal to the maximum value of $\sigma_t(\lambda)$ between time 0 and time τ . If this function satisfies $\sigma(\lambda_r) < 0$ for all r , then it follows that all $\sigma_t(\lambda_r) < 0$ also, and hence the system is stable.

Thus we are again able to define a master stability function that tells us whether the system is stable. Stability in this case means that small perturbations away from the synchronized state will die out and the system will remain synchronized. Once again the master stability function allows us to separate properties of the network from properties of the dynamics. The function depends only on the matrices $\alpha(t)$ and $\beta(t)$, which are defined by the dynamics through Eq. (17.103), while the eigenvalues λ_r depend only on the structure of the network. So for a given dynamics we can say what properties the network must have for stability or, conversely, for a given network we can say what properties the dynamics must have.

The downside of this approach is that it provides only a sufficient condition for stability and not a necessary one. Specifically, overall stability does not necessarily require that the system be instantaneously stable at every individual moment, as we have here assumed. It would be possible for perturbations around the synchronized state to grow at certain times, provided they decayed fast enough at others so that, on balance, they grow smaller over time. Thus, if the condition $\sigma(\lambda_r) < 0$ is satisfied for all r we can be certain that the synchronized state is stable, but if the condition is not satisfied it does not necessarily mean the synchronized state is unstable. In this respect our results here are weaker than those for the stability of fixed points.

It is possible to derive a full necessary and sufficient condition for stability, although in practice the calculations involved are often challenging. Briefly, one can define a function, or map, that takes as its argument the state of the system at time t , represented by the complete set of variables $\{\mathbf{x}^i(t)\}$ for all

nodes i and returns the state $\{\mathbf{x}^i(t + \tau)\}$ a time τ later, where τ is again the period of oscillation in the synchronized state. The periodic, synchronized solution $\mathbf{s}(t)$ is necessarily a fixed point of this function, since $\mathbf{s}(t + \tau) = \mathbf{s}(t)$, and, linearizing around this fixed point, one can then determine once more whether small perturbations will die away or grow. The condition for stability is given in terms of the eigenvalues of the linearized function [441].

The catch is that in practical situations it is rarely possible to calculate the map from t to $t + \tau$ exactly, since doing so would require us to solve the differential equations (17.96) describing the dynamics of the system, and if we were able to do that we probably would not be resorting to fixed-point techniques in the first place. In most cases, therefore, the best one can do is a numerical approximation. There is a substantial body of work, going by the name of *Floquet theory*, that offers ways to simplify the calculations and make them more robust and general, but ultimately one must still turn to numerics in the end, which limits the usefulness of the approach.

Many other details of network synchronization processes and many special cases have been studied. Another interesting area of investigation focuses on the question of whether a system will synchronize if the oscillators on the individual nodes are not identical. What happens if the oscillators have slightly different periods of oscillation, for instance? Normally we would then expect them not to be synchronized, to run at different rates, but it turns out that if the interaction between them is strong enough they will synchronize. There is a synchronization phase transition as we vary the strength of the interaction, between a weakly coupled, unsynchronized state and a strongly coupled, synchronized state. For a comprehensive discussion of this and a range of other interesting phenomena, the interested reader is encouraged to consult the review by Arenas *et al.* [30].

EXERCISES

17.1 Consider a dynamical system on a k -regular network (i.e., one in which every node has the same degree k) satisfying

$$\frac{dx_i}{dt} = f(x_i) + \sum_j A_{ij} g(x_i, x_j),$$

and in which the initial condition is uniform over nodes, so that $x_i(0) = x_0$ for all i .

a) Show that $x_i(t) = x(t)$ for all i where

$$\frac{dx}{dt} = f(x) + kg(x, x),$$

and hence that one need solve only one equation to solve the dynamics.

b) Show that for stability around a fixed point at $x_i = x^*$ for all i we require that

$$\frac{1}{k} > -\frac{1}{f'(x^*)} \left[\left(\frac{\partial}{\partial u} + \frac{\partial}{\partial v} \right) g(u, v) \right]_{u=v=x^*}.$$

17.2 Consider a dynamical system on an undirected network, with one variable per node obeying

$$\frac{dx_i}{dt} = f(x_i) + \sum_j A_{ij} [g(x_i) - g(x_j)],$$

as in Section 17.2.2. Suppose that the system has a symmetric fixed point at $x_i = x^*$ for all i .

a) Show, using results given in this chapter, that the fixed point is always stable if $f'(x^*) < 0$ and the largest degree k_{\max} in the network satisfies

$$\frac{1}{k_{\max}} > -2 \left[\frac{dg}{dx} / \frac{df}{dx} \right]_{x=x^*}.$$

b) Suppose that $f(x) = rx(1-x)$ and $g(x) = ax^2$, where r and a are positive constants.

Show that there are two symmetric fixed points for this system, but that one of them is always unstable.

c) Give a condition on the maximum degree in the network that will ensure the stability of the other fixed point.

17.3 The dynamical systems we have considered in this chapter have all been on undirected networks, but systems on directed networks are possible too. Consider a dynamical system on a directed network in which the sign of the interaction along an edge attached to a node depends on the direction of the edge, ingoing edges having positive sign and outgoing edges having negative sign. An example of such a system is a food web of predator–prey interactions, in which an ingoing edge indicates in-flow of energy to a predator from its prey and an outgoing edge indicates out-flow from prey to predator. Such a system can be represented by a dynamics of the form

$$\frac{dx_i}{dt} = f(x_i) + \sum_j (A_{ij} - A_{ji})g(x_i, x_j),$$

where g is a symmetric function of its arguments: $g(u, v) = g(v, u)$.

a) Consider a system of this form in which the in- and out-degrees of every node are equal to the same constant k . Show that such a system has a symmetric fixed point $x_i^* = x^*$ for all i satisfying $f(x^*) = 0$.

- b) Writing $x_i = x^* + \epsilon_i$, linearize around this fixed point to show that in the vicinity of the fixed point the vector $\epsilon = (\epsilon_1, \epsilon_2, \dots)$ satisfies

$$\frac{d\epsilon}{dt} = (\alpha \mathbf{I} + \beta \mathbf{M})\epsilon,$$

where $\mathbf{M} = \mathbf{A} - \mathbf{A}^T$. Determine the values of the constants α and β .

- c) Show that the matrix \mathbf{M} has the property $\mathbf{M}^T = -\mathbf{M}$. Matrices with this property are called *skew-symmetric*.
 d) If \mathbf{v} is a right eigenvector of a skew-symmetric matrix \mathbf{M} with eigenvalue μ , show that \mathbf{v}^T is a left eigenvector with eigenvalue $-\mu$. Hence by considering the equality

$$\mu = \frac{\mathbf{v}^\dagger \mu \mathbf{v}}{\mathbf{v}^\dagger \mathbf{v}} = \frac{\mathbf{v}^\dagger \mathbf{M} \mathbf{v}}{\mathbf{v}^\dagger \mathbf{v}},$$

where \mathbf{v}^\dagger is the Hermitian conjugate of \mathbf{v} , show that the complex conjugate of the eigenvalue is $\mu^* = -\mu$ and hence that all eigenvalues of a skew-symmetric matrix are imaginary.

- e) Show that the dynamical system is stable if $\text{Re}(\alpha + \beta \mu_r) < 0$ for all eigenvalues μ_r of the matrix \mathbf{M} , and hence that the condition for stability is simply $\alpha < 0$.

The last result means that the coupled dynamical system is stable at the symmetric fixed point if and only if the individual nodes are stable in the absence of interaction with other nodes.

17.4 The largest (most positive) eigenvalue κ_1 of the adjacency matrix of a k -regular graph, a Poisson random graph with mean degree c , and a star graph with n nodes, is k , $c + 1$, and $\sqrt{n - 1}$, respectively. Verify that the inequalities $\kappa_1 \geq \langle k \rangle$ and $\kappa_1 \geq \sqrt{\langle k^2 \rangle}$ from Eqs. (17.86) and (17.89) are satisfied in each of these cases.

17.5 Following the arguments of Section 17.2.2, the stability of a fixed point in certain dynamical systems on networks depends on the spectrum of eigenvalues of the adjacency matrix. Suppose we have a network that takes the form of an $L \times L$ square lattice, with each node labeled by its position vector $\mathbf{r} = (i, j)$ where $i, j = 1 \dots L$ are the row and column indices of the node. And suppose also that the system has periodic (toroidal) boundary conditions along its edges, such that the node at position $(i, 1)$ is adjacent to the node at (i, L) and the node at $(1, j)$ is adjacent to (L, j) .

- a) Consider the vector \mathbf{v} with one element for each node such that $v_{\mathbf{r}} = \exp(i\mathbf{k}^T \mathbf{r})$ for some vector \mathbf{k} . Show that \mathbf{v} is an eigenvector of the adjacency matrix provided

$$\mathbf{k} = \frac{2\pi}{L} \begin{pmatrix} n_1 \\ n_2 \end{pmatrix},$$

where n_1 and n_2 are integers.

- b) What range of values is permitted for the integers n_1 and n_2 ? Hence find the largest and smallest eigenvalues of the adjacency matrix.

17.6 Consider a network with an oscillator on every node. The state of the oscillator on node i is represented by a phase angle θ_i and the system is governed by dynamical equations of the form

$$\frac{d\theta_i}{dt} = \omega + \sum_j A_{ij} g(\theta_i - \theta_j),$$

where ω is a constant and the function $g(x)$ has $g(0) = 0$ and $g(x + 2\pi) = g(x)$ for all x .

- a) Show that the synchronized state $\theta_i = \theta^* = \omega t$ for all i is a solution of the dynamics.
- b) Consider a small perturbation away from the synchronized state $\theta_i = \theta^* + \epsilon_i$ and show that the vector $\epsilon = (\epsilon_1, \epsilon_2, \dots)$ satisfies

$$\frac{d\epsilon}{dt} = g'(0) \mathbf{L}\epsilon,$$

where \mathbf{L} is the graph Laplacian.

- c) Hence show that the synchronized state is stable against small perturbations if and only if $g'(0) < 0$.

CHAPTER 18

NETWORK SEARCH

A discussion of methods for searching networks for particular nodes or items, a process important for Web search and peer-to-peer networks, and for our understanding of the workings of social networks

IN CHAPTER 3 we saw a number of examples of networks that have information stored at their nodes: the World Wide Web, citation networks, peer-to-peer networks, and so forth. These networks can store large amounts of data but those data would be virtually useless without some way of searching through them for particular items. So important is it to be able to perform fast and accurate searches that the companies that provide the most popular search services are now some of the largest in their respective industries—Google, Thomson Reuters, LexisNexis—and constitute multibillion dollar international operations. In this chapter we examine some of the network issues involved in efficient searching and some implications of search ideas for the structure and behavior of networks.

18.1 WEB SEARCH

We have already discussed some aspects of how web search engines work in Sections 3.1 and 7.1.4. In this section we discuss the issue in more detail.

Traditional, or offline, web search is a multistage process. It involves first “crawling” the Web to find web pages and recording their contents, then creating an annotated index of those contents, including lists of words and estimates

of the importance of pages based on a variety of criteria. Then in the search process itself a user submits a text query to a search engine and the search engine extracts a list of pages matching that query from the index.

The process of web crawling by which web pages are discovered is interesting in itself and exploits the network structure of the Web directly. The crawler follows hyperlinks between web pages in a manner similar to the breadth-first search algorithm of Section 8.5. The basic process is described in Section 3.1. Practical web crawlers for big search operations employ many elaborations of this process, including:

- Searching in parallel at many locations on the Web simultaneously using many different computers;
- Placing the computers at distributed locations around the world to speed access times to pages coming from different places;
- Repeatedly crawling the same web pages at intervals of a few days or weeks to check for changes in page contents or pages that appear or disappear;
- Checking on pages more often if their contents have historically changed more often;
- Checking on pages more often if they are popular with users of the search engine;
- Heuristics to spot dynamically generated pages that can lead a crawler into an infinite loop or tree of pages and thereby waste time;
- Targeted crawling that probes more promising avenues in the network first;
- Altered behavior depending on requests from owners of specific sites, who often allow only certain crawlers to crawl their pages, or allow crawlers to crawl only certain pages, in order to reduce the load on their servers.

The processing of the raw crawler output also has interesting network-related elements. Early search engines simply compiled indexes of words or phrases occurring in web pages, so one could look up a word and get a list of pages containing it. Pages containing particular combinations of words could also be found by taking the sets of pages containing each individual word and forming the intersection of those sets. Indexes can be extended by adding annotations indicating, for example, how often a word appears on a page or whether it appears in the page title or in a section heading, which might indicate a stronger connection between that word and the subject matter of the page. Such annotations allow the search engine to make choices about which are the pages most relevant to a given query. Even so, search engines based solely on

indexes and textual criteria of this sort do not return very good results and have been superseded by more sophisticated technology.

Modern search engines do still use indexes in their search process, but only as a first step. A typical modern search engine will use an index to find a set of candidate pages that might be relevant to a given query and then narrow that set down using other criteria, some of which may be network-based. The initial set is usually chosen deliberately to be quite broad. It will typically include pages on which the words of the query appear, but also pages on which they don't appear but that link to, or are linked to by, pages that do contain the query words. The net result is a set of pages that probably includes most of those that might be of interest to the user submitting the query, but also many irrelevant pages as well. The strength of the search engine, its ability to produce useful results, therefore rests primarily on its ability to narrow the search within this broad set.

The classic example of a criterion used for narrowing web searches comes from the Google search engine, which makes use of the eigenvector centrality measure known as PageRank, discussed in Section 7.1.4. PageRank accords pages a high score if they receive hyperlinks from many other pages, but does so in a way such that the credit received for a link is higher if it comes from a page that is itself highly ranked. PageRank, however, is only one of many elements that go into the formula Google uses to rank web pages. Others include traditional measures such as frequency of occurrence of query words in the page text and position of occurrence (near the top or bottom, in titles and headings, etc.), as well as occurrence of query words in "anchor text" (the highlighted text that denotes a hyperlink in a referring page) and previous user interest in a particular page (whether people selected this page from the list of search results on other occasions when the same text query, or a similar one, was entered).

Google gives each web page in the initial set a score that is a weighted combination of these elements and others. The particular formula used is a closely guarded secret and is moreover constantly changing, partly just to try and improve results, but also to confound the efforts of web page creators, who try to increase their pages' ranking by working out what particular elements carry high weight in Google's formula and incorporating those elements into their pages.

An important point to appreciate is that some parts of the score a page receives depend on the particular search query entered by the user—frequency of occurrence of query words, for instance—but others, such as PageRank, do not. This allows Google's computers (or their counterparts at other search companies) to calculate the latter parts "offline," meaning they are calculated

ahead of time and not at the time of the query itself. This has some advantages. PageRank, for instance, is computationally intensive to calculate and it saves a lot of time if you only have to calculate it once. But there are disadvantages too. PageRank measures the extent to which people link to a given web page, but people may link to a page for many reasons. Thus a page may have a high PageRank for a reason unrelated to the current search query. A page whose text makes mention of two or more different topics (and many do) may be a crucial authority on one topic but irrelevant on another, and PageRank cannot distinguish between the two.

One could imagine a version of PageRank that was specific to each individual query. One could calculate a PageRank score within just the subnetwork formed by the set of pages initially selected from the index to match the query. But this would be computationally expensive and it's not what Google does. As a result it is not uncommon for a page to be ranked highly in a particular search even though a casual human observer could quickly see that it was irrelevant to the search topic. In fact, a large fraction of "bad" search results returned by search engines probably fall in this category: they are pages that are important in some context, but not in the context of the specific search conducted.

The overall process behind searches on Google and similar large search engines is thus as follows [82]. First the Web is crawled to find web pages. The text of those web pages is processed to create an annotated index, and the link structure of the hyperlinks between them is used to calculate a centrality score or scores for each page, such as PageRank in Google's case or (presumably) some similar measure for other search engines. When a user enters a query the search engine extracts a deliberately broad set of matching pages from the index, scores them according to various query-specific measures such as frequency of occurrence of the query words, then combines those scores with the pre-computed centrality measure and possibly other pre-computed quantities, to give each page in the set an overall score. Then the pages are sorted in order of their scores and those with the highest scores are transmitted to the user. Typically, only a small number of the highest-scoring pages are transmitted—say the first ten—but with an option to see lower-scoring pages if necessary.

Despite the reservations mentioned above, this system works well in practice, far better than early web search engines based on textual content alone, and provides useful search results for millions of computer users every day.

18.2 SEARCHING DISTRIBUTED DATABASES

Some information networks take the form of distributed databases. A typical example is a peer-to-peer file-sharing network, in which individual computers

in the network each store a subset of the data stored in the network as a whole. The form and function of peer-to-peer networks were described in Section 3.3.1.

The “network” in a peer-to-peer network is a virtual one, in which individual computers maintain contacts with a subset of others, not necessarily those with which they have direct physical data connections. In this respect peer-to-peer networks are somewhat similar to the World Wide Web, in which the hyperlinks between websites are virtual links chosen by a page’s creator and their topology need have nothing to do with the topology of the underlying physical Internet. Indeed, the World Wide Web is itself, in a sense, a distributed database, storing information in the pages at its nodes, but web search works in a fundamentally different way from search in other distributed databases, so we treat the two separately.

Search is a fundamental problem in peer-to-peer networks and other distributed databases: without a way to find specific items among those stored at the many nodes of the network, the system is essentially useless. One solution is to copy the web search approach of Section 18.1 and construct a central index of all items and then search that index for items of interest. As discussed in Section 3.3.1, however, most peer-to-peer networks don’t go this route, but instead make use of distributed search techniques in which the search task is shared among the computers in the network via messages passed along network edges. Indeed the performance of such distributed searches is the primary reason for linking the nodes into a network in the first place and there are some interesting connections between the structure of the network and the efficiency with which searches can be performed.

Suppose that we have a peer-to-peer network composed of n individual computers and each computer is linked by virtual connections to a selection of the others, where “linked” in this context merely means that these others are the ones with which a computer has agreed to communicate directly in the course of performing searches. There is no reason in principle why a computer could not communicate with *all* others if it wanted to, but in practice this is usually too demanding, and limiting the number of network neighbors a computer has makes the task more tractable.

The simplest form of distributed search, used in some of the earliest peer-to-peer networks, is a version of the breadth-first search algorithm described in Section 8.5 (where it was used for finding network components and shortest paths). Under this approach, a user gives the computer a search term, such as the name of a computer file, and the computer sends a query to each of its neighbors in the peer-to-peer network, asking if they have the file in question. If they do, they send the file to the computer that made the query and the search is complete. If they don’t, then they send a further query to each of *their*

neighbors asking for the file. Any neighbor that has seen the query before, such as the computer that originated it in the first place, ignores it. All others check to see whether they have the requested file and send it back to the originating computer if they do. If not, they pass the query on to their neighbors, and so on.

This simple strategy certainly works and it has some advantages. For instance, assuming that the network displays the small-world effect, the number of steps we will have to take in our breadth-first search will be small even when the network is large (typically increasing only logarithmically with n —see Section 11.7). This means that most searches will take only a short amount of time to find the desired file.

But there are some serious disadvantages with the approach as well. First, as we have described it the search doesn't actually stop when the target file is found. There is no mechanism to inform computers that the file has been found and that they don't need to pass the query on to anyone further. This problem can be fixed relatively easily, however, for example by requiring each computer receiving the query to check with the originating computer to see whether the file has been found before they do anything else.

A more serious problem is that the messages transmitted in the process of spreading a query across the network quickly add up to a huge amount of data and can easily overwhelm the capacity of the computers involved. Assuming a worst-case scenario in which a desired file exists on only a single computer in the network, we will, on average, have to pass our query to half of all computers before we find the file. That means the number of messages sent in the course of one query is $O(n)$. If the average user performs queries at some constant rate r per unit time, then the overall rate of queries for all n users is $rn = O(n)$. Thus, the total number of messages sent per unit time is $O(n) \times O(n) = O(n^2)$ and the number of messages sent per computer per unit time is, on average, $O(n^2)/n = O(n)$, which goes up linearly with the size of the network. This means that, no matter how much bandwidth our computers have to send and receive data, it will in the end always become swamped if the network becomes large enough. And peer-to-peer networks can become extremely large. Some of the largest have millions of users.

Luckily, this worst-case scenario is not usually realized. It is, in fact, rarely the case that an item of interest exists on only one computer in a network. In a typical peer-to-peer network most items exist in many places. Indeed, assuming that some fraction of the user population likes or needs each item, it is more reasonable to suppose that any given item appears on some fixed fraction c of the nodes in the network, so that the total number of copies cn goes up as the size of the network increases. If this is the case, and assuming

See Section 10.2 for a discussion of the small-world effect.

for the moment that the value of c is the same for every item, then one will have to search on average only $1/c$ nodes before finding a copy of an item. This means that the total number of query messages sent over the network per unit time is $O(n/c)$ and the number per computer per unit time is $O(1/c)$, which is just a constant and does not increase with increasing network size.

A more realistic calculation allows for the fact that some items are more popular than others. Suppose that the factors c , which are proportional to popularity, have a distribution $p(c)$, meaning that the probability of falling in the interval c to $c + dc$ is $p(c) dc$. Also important to note is that not all items are searched for with equal frequency. Indeed a more reasonable assumption is that they are searched for with frequency proportional to their popularity, meaning that the probability of a query asking for an item with popularity c is proportional to c and the probability of asking for any one of the items with popularity in the interval c to $c + dc$ is proportional to $cp(c) dc$. We can calculate the constant of proportionality by noting that this probability must be normalized to one over all values of c , so the correctly normalized distribution must be

$$\frac{cp(c) dc}{\int cp(c) dc} = \frac{cp(c) dc}{\langle c \rangle}, \quad (18.1)$$

where $\langle c \rangle = \int cp(c) dc$ is the mean value of c .

Then the average number of nodes we have to examine before we find the item corresponding to a typical query is

$$\int_0^1 \frac{1}{c} \frac{cp(c) dc}{\langle c \rangle} = \frac{1}{\langle c \rangle}, \quad (18.2)$$

where we have made use of the fact that $\int p(c) dc = 1$. Hence the number of query messages sent or received per computer per unit time is $O(1/\langle c \rangle)$, which is again a constant as network size becomes large.

In principle, therefore, if a node can handle messages at the rate given by Eq. (18.2), then the network should go on functioning just fine as its size becomes large. In practice, however, there can still be problems. The main difficulty is that nodes in the network vary enormously in their bandwidth capabilities. Most nodes have relatively slow communications with the network, i.e., low bandwidth, while a few have much better, higher-bandwidth connections. This means that even if bandwidth requirements per node are reduced to a constant as above, the network will still run at a speed dictated by the majority slow nodes, making queries slow and possibly overwhelming the capacity of some nodes.

To get around this problem, most modern peer-to-peer networks make use of *supernodes* (also called *superpeers*). Supernodes are high-bandwidth nodes

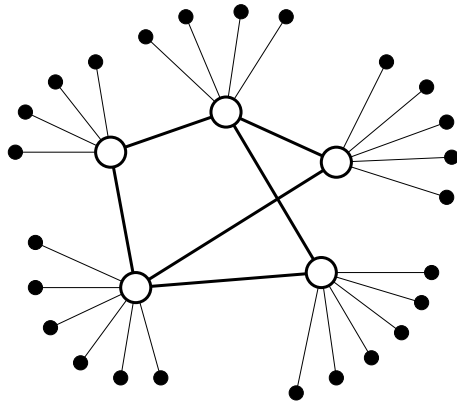


Figure 18.1: The structure of a peer-to-peer network with supernodes. Client nodes (filled circles) are connected to a network of supernodes (open circles) that have above-average network bandwidth and hence can conduct searches quickly.

chosen from the larger population in the network and connected to one another to form a supernode network over which searches can be performed quickly—see Fig. 18.1.

A supernode acts a little like a local exchange in a telephone network (see Section 2.2). Each normal user, or *client*, in the network attaches to a supernode (or sometimes to more than one) that acts as their link to the rest of the network. Each supernode has a number of such clients and the clients communicate to the supernode a list of the files or other data items they possess so that the supernode can respond appropriately to search queries on their behalf. An individual client wanting to perform a search then sends a search query to the local supernode, which initiates a breadth-first search of the network of supernodes to find the desired item. Since the supernodes possess records of all the items that the clients have, the entire search can be performed on the network of supernodes alone and no client resources are used at all. And since the supernodes are deliberately selected to have fast network connections, the search runs at the speed of the quickest nodes in the network.

In practice schemes like this work quite well—well enough to be in wide use in peer-to-peer networks of millions of users. More sophisticated schemes have been devised that in theory could work better still—an example is the “Chord” system proposed by Stoica *et al.* [438]—but such systems have yet to find widespread adoption since the more traditional supernode approach appears to work well enough for practical purposes.

18.3 SENDING MESSAGES

A different variation of the distributed search problem is the problem of getting a message to a particular node in a network. The classic example of this problem is Stanley Milgram's "small-world" experiment, described in Section 4.6. In this experiment participants were asked to get a message to a specific target individual by passing it from acquaintance to acquaintance through the social network. Milgram famously found that messages that reached the target took only about six steps to do so, which is the origin of the popular concept of the "six degrees of separation." As mentioned in Section 4.6, however, there is another perhaps more surprising implication of the experiment, first pointed out by Kleinberg [266], which is that short paths not only exist in the network but that people are remarkably good at finding them. Of course if one knows the structure of an entire network then one can find short paths directly using, for example, the breadth-first search method of Section 8.5.5. Participants in Milgram's experiment, however, did not know the whole network and probably only knew a very small part of it, and yet they were still able to get a message rapidly to the desired target.

This observation raises a number of interesting questions. How, in practice, did people find these short paths to the target? Can we come up with an algorithm that will do the job efficiently? How does the performance of that algorithm depend on the structure of the network? In the following sections we consider two different models of the message passing process that address these questions. As we will see, these models suggest that social (or other) networks must have a very particular type of structure if one wants to be able to find short paths easily without global knowledge of the network.

18.3.1 KLEINBERG'S MODEL

The instructions to the participants in Milgram's experiment were that upon receiving the message (actually a small booklet or "passport" sent through the mail), they were to forward it to an acquaintance whom they believed to be closer to the target than they were. The definition of "closer" was left vague, however, and one of the first things we need to do if we want to model the mechanics of the experiment is to decide on a practical definition.

An illuminating attempt at modeling Milgram's experiment was made by Kleinberg [266, 267], who employed a variant of the small-world model of Section 12.11.8, as shown in Fig. 18.2. As in the standard small-world model, it has a ring of nodes around the edge plus a number of "shortcut" edges that connect node pairs at random points around the ring. In Kleinberg's model all

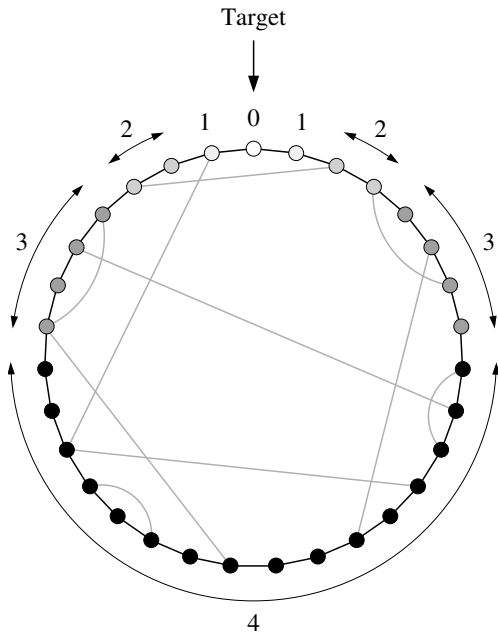


Figure 18.2: The variant of the small-world model used to model message passing. Nodes are connected around a ring and shortcuts added between them as in the normal small-world model. However, the shortcuts are now biased so that there are more of them connecting nearby nodes than distant nodes. The strength of the bias is controlled by the parameter α . In the proof given in the text, the nodes are divided into numbered classes, class 0 consisting of just the target node and higher classes radiating out from the target, each successive class containing twice as many nodes as the previous one.

nodes are connected to their two immediate neighbors around the ring ($c = 2$ in the notation of Section 12.11.8) and Kleinberg made use of the connections in the ring to define the “closeness” of nodes for the purposes of the message-passing experiment. He proposed that individuals in the network, represented by nodes, are aware of the distance around the ring to other individuals, and hence can say when one of their acquaintances is closer to the target node than they are in this sense. One could imagine, for instance, that the ring represents geographic space and distance around it is a measure of how close, geographically speaking, an individual is to the target.

In his calculations Kleinberg considered a *greedy algorithm* for message passing in which each individual receiving a message passes it on to the one of their neighbors who is closest to the target in the sense above. This algorithm is guaranteed always to get the message to the target eventually. Every individual has at least one neighbor who is closer to the target in the Kleinberg sense than they are—their neighbor around the ring in the direction towards the target. Thus on each step of the message-passing process the message is guaranteed to get at least one step closer to the target and hence it must eventually get to the target. In the worst case, individuals simply pass the message around the ring until it reaches its destination, but generally we can expect to do better than this because of the shortcuts. The question is how much better. Kleinberg showed

that it is possible for the greedy algorithm to find the target node in $O(\log^2 n)$ steps, but that it can do so only for particular choices of the arrangement of the shortcuts.

Kleinberg considered a one-parameter family of models that generalizes the standard small-world model by allowing for different arrangements of the shortcuts.¹ Instead of assuming that shortcuts are placed uniformly at random, we assume (not unreasonably) that people have more acquaintances among those close to them (in the sense defined above) than among those far away. By analogy with the standard small-world model let us place a total number of shortcuts equal to p times the number of edges in the ring itself, which in this case is just n . Since each shortcut has two ends this means that the average number of shortcuts attached to each node will be $2p$ (and the actual number will be Poisson distributed with mean $2p$). Where we differ from the standard small-world model is in the lengths of the shortcuts. Shortcuts are still located at random around the ring, but they are chosen so that the probability of a particular shortcut covering a distance r is $Kr^{-\alpha}$, where α is a non-negative constant and K is a normalizing constant. That is, for each shortcut we first choose its length r from this distribution, then we place the shortcut, spanning a distance exactly r , at a position chosen uniformly at random around the ring. If $\alpha = 0$ then we recover the standard small-world model of Section 12.11.8, but more generally, for $\alpha > 0$, the model has a preference for making connections between nearby nodes rather than distant ones.

Note that the probability that a particular shortcut connects a specific pair of nodes a distance r apart is equal to $Kr^{-\alpha}/n$, which is the probability $Kr^{-\alpha}$ that the shortcut has length r multiplied by the probability $1/n$ that out of the n possible positions around the ring it happens to fall in the one that connects the two nodes in question. Given that there are np shortcuts in the whole network, this means that the total probability of having a shortcut between a given pair of nodes is $np \times Kr^{-\alpha}/n = pKr^{-\alpha}$. (More correctly, this is the expected number of such shortcuts, but so long as the number is small it can be interpreted as a probability to an excellent approximation.)

The normalizing constant K is fixed by the condition that every shortcut

¹The model we use is a somewhat simplified version of Kleinberg's. His model used a two-dimensional lattice instead of a one-dimensional ring as the underlying structure for the network, though the calculations work in the same way in either case. Our model also places both ends of each shortcut at random, where Kleinberg's fixed the number of shortcuts attached to each node to be constant and also made them directed rather than undirected.

must have *some* length, and that all lengths lie between 1 and $\frac{1}{2}(n-1)$, so that²

$$K \sum_{r=1}^{\frac{1}{2}(n-1)} r^{-\alpha} = 1. \quad (18.3)$$

We can approximate the sum by an integral using the trapezoidal rule of Eq. (13.103) thus:

$$\begin{aligned} \sum_{r=1}^{\frac{1}{2}(n-1)} r^{-\alpha} &\simeq \int_1^{\frac{1}{2}(n-1)} r^{-\alpha} dr + \frac{1}{2} + \frac{1}{2} \left[\frac{1}{2}(n-1) \right]^{-\alpha} \\ &= \frac{\left[\frac{1}{2}(n-1) \right]^{1-\alpha} - 1}{1-\alpha} + \frac{1}{2} + \frac{1}{2} \left[\frac{1}{2}(n-1) \right]^{-\alpha}, \end{aligned} \quad (18.4)$$

which gives

$$K \simeq \begin{cases} (1-\alpha)\left(\frac{1}{2}n\right)^{\alpha-1} & \text{for } \alpha < 1, \\ 1/\ln \frac{1}{2}n & \text{for } \alpha = 1, \\ 2(\alpha-1)/(\alpha+1) & \text{for } \alpha > 1, \end{cases} \quad (18.5)$$

as n becomes large.³

We can now show that, for suitable choice of α , the greedy algorithm on this network can indeed find a given target node quickly. The proof is as follows. Suppose, without loss of generality, that the target node is at the top of the ring, as depicted in Fig. 18.2, and let us divide up the other nodes into classes according to their distance from the target. Class 0 consists of just the target itself. Class 1 consists of all nodes distance $d = 1$ from the target around the ring, of which there are two. Class 2 consists of nodes with distances in the range $2 \leq d < 4$, class 3 of nodes $4 \leq d < 8$, and so forth. Each class is double the size of the previous one. In general, class k consists of nodes at distances $2^{k-1} \leq d < 2^k$ and contains $n_k = 2^k$ nodes. (For simplicity, let us assume that the total number n of nodes is a power of two, minus one, so that everything works out neatly.)

Now consider a message being passed through the network according to our greedy algorithm and suppose that at a particular step of the process the message is at a node of class k . How many more steps will it take before the message leaves class k and passes into a lower class? The total number of nodes

²The maximum length of a shortcut is $\frac{1}{2}(n-1)$ if n is odd and $\frac{1}{2}n$ if n is even. We will assume that n is odd in this case, which avoids some small annoyances in the derivations.

³Note that both the numerator and denominator of the fraction in Eq. (18.4) vanish at $\alpha = 1$, so one must use l'Hopital's rule to extract the limiting value. The same goes for Eq. (18.9).

in lower classes is

$$\sum_{m=0}^{k-1} n_m = \sum_{m=0}^{k-1} 2^m = 2^k - 1 > 2^{k-1}, \quad (18.6)$$

and from Fig. 18.2 we can see that all of these nodes are, at most, a distance $3 \times 2^{k-1} - 2 < 2^{k+1}$ from the node in class k that currently holds the message. Thus, the probability of the node with the message having a shortcut to a particular one of these nodes in lower classes is at least $pK 2^{-(k+1)\alpha}$, and the probability of having a shortcut to any of them is at least $pK 2^{k-1-(k+1)\alpha}$.

If our node has no shortcut that takes the message out of class k , then, in the worst case, it simply passes the message to another node in class k that is closer to the target, either via a shortcut or by passing around the ring. Using the probability above, the expected number of such moves made before we find a shortcut that takes us out of class k is at most

$$\frac{1}{pK 2^{k-1-(k+1)\alpha}} = \frac{1}{pK} 2^{\alpha+1} 2^{(\alpha-1)k}. \quad (18.7)$$

Finally, again in the worst case, the message will pass through each of the classes before reaching the target. Excluding the target itself, the classes run from $k = 1$ to $\log_2(n+1) - 1$ and summing over them we find that an upper bound on the expected number of steps ℓ needed to reach the target is

$$\begin{aligned} \ell &\leq \frac{1}{pK} 2^{\alpha+1} \sum_{k=1}^{\log_2(n+1)-1} 2^{(\alpha-1)k} = \frac{1}{pK} 2^{\alpha+1} \frac{2^{(\alpha-1)\log_2(n+1)} - 2^{\alpha-1}}{2^{\alpha-1} - 1} \\ &= \frac{1}{pK} 2^{\alpha+1} \frac{(n+1)^{\alpha-1} - 2^{\alpha-1}}{2^{\alpha-1} - 1}. \end{aligned} \quad (18.8)$$

Making use of Eq. (18.5) for the constant K and taking the limit of large n we find that asymptotically

$$\ell \leq \begin{cases} A n^{1-\alpha} & \text{if } \alpha < 1, \\ B \log^2 n & \text{if } \alpha = 1, \\ C n^{\alpha-1} & \text{if } \alpha > 1, \end{cases} \quad (18.9)$$

where A , B , and C are constants depending on α and p , but not n , whose rather complicated values we can work out from Eqs. (18.5) and (18.8) if we want to.

Since Eq. (18.9) gives an upper bound on ℓ , this result guarantees that for the particular case $\alpha = 1$ we will be able to find the target node in a time that increases at most as $\log^2 n$ with the size of the network. This is not quite as good as $\log n$, which is the actual length of the shortest path in a typical network, but it is still a slowly growing function of n and it would be fair to claim that the

small-world experiment would succeed in finding short paths in a network that had $\alpha = 1$. Thus it is possible, provided the network has the correct structure, for a simple strategy like the greedy algorithm, in which nodes have knowledge only of their immediate network neighborhood, to produce results similar to those observed by Milgram in his experiment.

On the other hand, if $\alpha \neq 1$ then Eq. (18.9) increases as a power of n , suggesting that it would take much longer in such networks to find the target node. In particular, for the original small-world model of Section 12.11.8, which corresponds to $\alpha = 0$, Eq. (18.9) grows linearly with n , suggesting that the Milgram experiment could take millions of steps to find a target in a social network of millions of people. Equation (18.9) is only an upper bound on the time taken, so if one is lucky one may be able to find the target faster. For instance, if the message starts at a node that happens to have a shortcut directly to the target node then one can find the target in a single step. However, Kleinberg [267] was also able to prove that the *average* time it takes to find the target increases at least as fast as a power of n except in the special case $\alpha = 1$, so in general the greedy algorithm for $\alpha \neq 1$ will not work well.⁴

These results tell us two things. First, they tell us that it is indeed possible for the small-world experiment to work as observed even if the participants don't know the details of the whole network. But, second, they tell us that, at least within the context of the admittedly not very realistic model used here, the experiment only works for certain special values of the parameters of the network. Thus, the success of Milgram's experiment suggests not only that, as Milgram concluded, there are short paths in social networks, but also that the networks have a particular structure that makes path finding possible.

18.3.2 A HIERARCHICAL MODEL FOR MESSAGES

While interesting, the results of the previous section are not wholly convincing because the model is clearly not a realistic one. People don't live around a circle with just a few shortcuts to others, and message passing doesn't work because people know where others live on the circle.

So can we derive similar results for a more realistic network model? To answer this question let us first ask how the transmission of messages *does* work. We can get a hint from the "reverse small-world" experiments of Killworth

⁴In fact, since Kleinberg was studying a two-dimensional version of the small-world model, his result was for $\alpha = 2$, not $\alpha = 1$. In general, on a small-world network built on a d -dimensional lattice, the greedy algorithm succeeds in finding the target in time $O(\log^2 n)$ only when $\alpha = d$ and for all other values takes time increasing at least as a power of n .

and Bernard [57,260] discussed in Section 4.6. Recall that in these experiments researchers asked subjects to imagine that they were participating in Milgram's small-world experiment and then asked them what information they would want to know about the target in order to make a decision about whom to pass their message to. Killworth and Bernard found that three pieces of information were sought more often than any others, and by almost all subjects: the name, occupation, and geographic location of the target.

The target's name is an obvious requirement in the small-world experiment, since it's needed to recognize the target when you find him or her. Beyond that, however, it probably doesn't play much role in the message passing, except perhaps in cultures where names can give a clue to the location or social status of an individual. Occupation and geographic location, on the other hand, are of great use in deciding how to forward a message, and these appear to be the primary pieces of information participants in the experiment use.

Take geographic location as an example. How would one use information on geography to route a message? Presumably, one would attempt to pass the message to someone closer geographically to the target than oneself. Suppose, for instance, that the target lives, as Milgram's did, in a suburb of the city of Boston, Massachusetts, in the United States. A participant in Britain, attempting to get a message to this target, might perhaps forward it to someone they knew in the US, say in New York. That person might forward it in turn to someone they knew in the state of Massachusetts, who would forward it to someone in Boston, who would forward it to the target's specific suburb, and so forth. At each step in the process, the participants narrow the search to a smaller geographic area until, with luck, the area is so small that someone there knows the target individual directly.

In a sense, this is what happens in Kleinberg's model. In Section 18.3.1 we divided Kleinberg's circle into zones or classes that get ever smaller as they close in on the target and showed that under suitable circumstances it takes only a small number of steps of the message-passing process to find a connection from one class to the next smaller one. Since the number of classes is logarithmic in the size of the network, this means that it also takes only a small number of steps overall to home in on the target. Kleinberg's network structure was unrealistic, but the basic idea of progressively narrowing the field is a good one and we would like to find a more realistic network model to which the same type of argument can be applied.

Such a model is the hierarchical model of Watts *et al.* [465], in which the interplay of social structure and geographic or other dimensions is represented by a tree or dendrogram that represents the progressive division of the par-

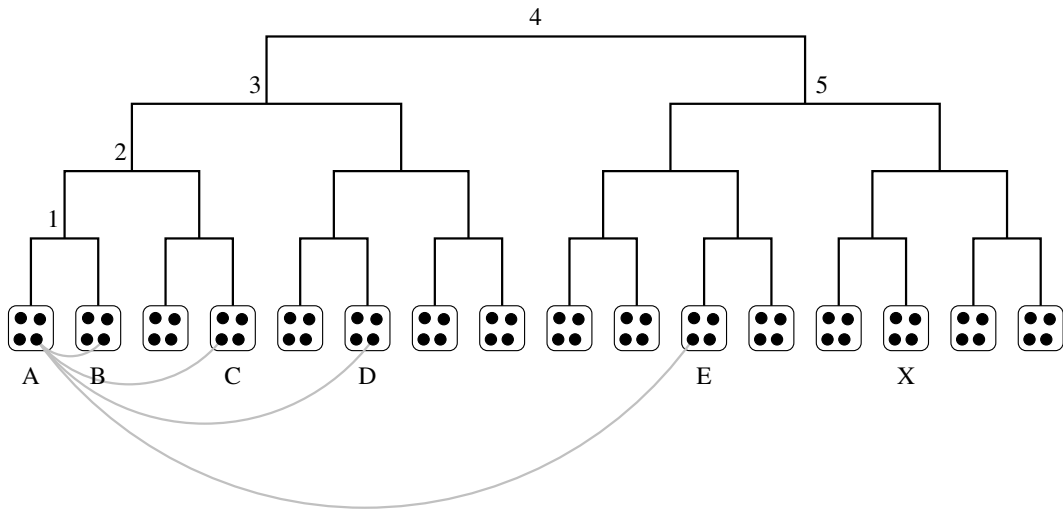


Figure 18.3: The hierarchical model of Watts *et al.* Small groups of individuals (boxes) are divided up in a hierarchical structure represented by a binary tree, which might, for instance, correspond to the hierarchical division of geographic space into countries, regions, towns, and so forth. The hierarchy dictates which social connections (indicated by curves) are most likely. A node in group A, for instance, is most likely to be connected to others close to it in the tree (B, C) and less likely to be connected to those further away (D, E).

ticipants into smaller and smaller groups.⁵ In the context of geography, for example, the world might be divided into countries, the countries into regions, the regions into cities and towns, and so forth. The division ends when we reach units so small that it can reasonably be assumed that everyone knows everyone else—a single family, for instance.

The divisions can be represented as shown in Fig. 18.3. The tree in this case is a binary tree. Each branch splits in two, then in two again, and so forth. In the real world branches might split into more than two parts. There are more than two countries in the world after all. However, the binary tree is the simplest case to study (and the one studied by Watts *et al.*), and the analysis given here for the binary case can be generalized to other cases quite easily.

Let us also assume that the groups at the bottom of the tree all have the same size g . Again this is a simplification, but a useful one that does not have a major effect on the results. If the total number of individuals in the network is n then the number of groups is n/g , and the number of levels in the tree is $\log_2(n/g)$.

⁵A similar model was also proposed independently by Kleinberg [268].

The model makes two other important assumptions. First, it assumes that people measure distance to a target individual in terms of the tree, and more specifically in terms of the lowest common ancestor in the tree that they share with the target. That is, people are able to tell when someone lives in the same country as themselves, or the same region or town, but do not have any detailed information beyond that. This is a more conservative assumption than is made by Kleinberg's model. In Kleinberg's model it is assumed that people know their exact geometric distance to the target, no matter where in the world the target lives. In the present model people have a more coarse-grained impression of how close they are to the target.

The second assumption in the model of Watts *et al.* is that the social network itself is correlated with the hierarchical tree structure so that people who are closer together in the tree, in the sense of sharing a lower common ancestor, are also more likely to be acquainted. (The model is similar in this respect to the hierarchical random graph discussed in Section 14.7.2.) Thus people are more likely to know others in their own country than in other countries, more likely to know others in their own town than in other towns, and so forth. A few sample acquaintances are represented by the curves at the bottom of the figure.

Thus there are really two networks present in this model. There is the "real" network of actual acquaintances represented by the curves, and a "shadow" network, the hierarchical tree, which is not a network of acquaintances but which influences the acquaintance network and of which individuals are somewhat aware, in the sense that they know how close they are to others in the tree.

An important point to note about this model is that although an individual is less likely to know others far away in the tree, there are also more such far-away individuals than there are ones close by, and the two effects cancel out to some extent so that it is quite possible for a given individual to know others who are both near and far. The people who live on your street, for instance, are close by, so you are likely to know them, but they are also few in number. By contrast, other countries may be far away from you, but they contain a lot of people, so even though you are not very likely to know any particular inhabitant, it is nonetheless quite possible that you know at least one out of the whole population. This behavior is crucial to making the message-passing experiment work.

Consider an individual in group A in Fig. 18.3. Let us suppose that, because of the effect described above, this individual has at least one acquaintance at every "distance" in the tree, i.e., one acquaintance in every subtree of the hierarchy with which they share a common ancestor. That is, they know one of the individuals in group B, the group with which they share ancestor 1,

and they also know someone in one of the two groups with whom they share ancestor 2—say someone in group C—and so on through groups D and E as shown. And suppose also that a similar pattern of acquaintances holds for every other individual in the network: everyone knows at least one person in every subtree with whom they share a common ancestor.

Now consider a greedy algorithm for message passing on this network. Suppose the message starts at a node in group A and, as before, the holder of the message at each step passes it to an acquaintance closer to the target than they are, distance now being measured in the sense of the hierarchical tree as described above. Suppose, for instance, that the target node is in group X, which shares a common ancestor with A only at the highest and coarsest level marked 4 in the figure. That is, the target is in the opposite subtree of ancestor 4 from A. By hypothesis, the individual holding the message knows this and hence knows that in order to get the message closer to the target they must pass it to someone in that opposite subtree. Luckily, under the assumption above they always have such an acquaintance, in this case in group E. So they pass the message to their friend in group E. The friend now notes that the target X is in the subtree with whom they share the common ancestor marked 5 and hence knows that they must pass the message to a neighbor in that subtree to get it closer to the target. Again, by definition, they have at least one such neighbor, so they pass the message to that neighbor. And so the process proceeds. At each step we narrow down our search to a smaller subtree of the overall network, or equivalently we move to a lower level in the hierarchy, pivoting about a common ancestor. But the total number of levels in the hierarchy is $\log_2(n/g)$ and hence this is the maximum number of steps that the process will take to reach the target. In this model, therefore, the message always reaches its target in a logarithmic number of steps.

It is not, however, very realistic to assume that each individual in the network knows at least one person at each distance. Watts *et al.* considered a more realistic probabilistic model in which there is a probability p_m of two individuals knowing one another when their lowest common ancestor is at level m in the tree. The levels are numbered so that $m = 1$ for groups that are immediately adjacent, as A and B are in Fig. 18.3, and increase by one for each higher level up to a maximum of $m = \log_2(n/g)$ at the top of the tree.

Watts *et al.* considered the particular choice

$$p_m = C 2^{-\beta m}, \quad (18.10)$$

where C and β are constants.⁶ So long as β is positive this choice gives, as

⁶Watts *et al.* actually wrote the expression as $Ce^{-\beta m}$, but the difference is only in the value of β

desired, a lower probability of acquaintance with more distant individuals, the exact rate of variation being controlled by the value of β . The parameter C controls the overall number of acquaintances that each individual has.

The number of nodes with which any given node shares its ancestor at level m is $2^{m-1}g$ and hence the expected number of such nodes that it will be connected to is

$$2^{m-1}gp_m = \frac{1}{2}Cg 2^{(1-\beta)m}. \quad (18.11)$$

Summing over all levels, the total expected number of acquaintances an individual has, their average degree in the network, we find to be

$$\langle k \rangle = \frac{1}{2}Cg \sum_{m=1}^{\log_2(n/g)} 2^{(1-\beta)m} = \frac{1}{2}Cg \frac{2^{(1-\beta)\log_2(n/g)} - 1}{1 - 2^{\beta-1}} = \frac{1}{2}Cg \frac{(n/g)^{1-\beta} - 1}{1 - 2^{\beta-1}}. \quad (18.12)$$

Thus the constant C is given by

$$C = \frac{2\langle k \rangle}{g} \frac{1 - 2^{\beta-1}}{(n/g)^{1-\beta} - 1}. \quad (18.13)$$

In the limit of large n this simplifies to

$$C = \begin{cases} (2\langle k \rangle/g)(1 - 2^{\beta-1})(n/g)^{\beta-1} & \text{for } \beta < 1, \\ (2\langle k \rangle/g)/\log_2(n/g) & \text{for } \beta = 1, \\ (2\langle k \rangle/g)(2^{\beta-1} - 1) & \text{for } \beta > 1. \end{cases} \quad (18.14)$$

If a particular node receives a message and wants to pass it to the opposite subtree at level m , it can do so provided it has an acquaintance in that subtree. If (18.11) is small, however, then most likely it will not, in which case the best it can do is to pass the message to someone else in the subtree it is already in, who can then repeat the process. The expected number of times this will happen before the message finds a person who does have a neighbor in the opposite subtree is given by the reciprocal of (18.11), which is $(2/Cg)2^{(\beta-1)m}$. Summing this over all levels, we find that the total expected number of steps to reach the target is

$$\ell = \frac{2}{Cg} \sum_{m=1}^{\log_2(n/g)} 2^{(\beta-1)m} = \frac{2}{Cg} \frac{2^{(\beta-1)\log_2(n/g)} - 1}{1 - 2^{1-\beta}} = \frac{2}{Cg} \frac{(n/g)^{\beta-1} - 1}{1 - 2^{1-\beta}}. \quad (18.15)$$

and we find the definition (18.10) to be more convenient.

It is possible that the node holding the message will have a neighbor neither in the opposite subtree nor in its own subtree. If this happens then the node has only neighbors further from the target than it is and none nearer. In this case the Milgram experiment fails—recall that participants were asked to pass the message to someone closer to the target. This, however, is not necessarily unrealistic. As Watts *et al.* point out, this presumably does happen in the real experiment sometimes, and moreover it is well documented that many messages, a majority in fact, got lost and never reached their target. For messages that do get through, however, Eq. (18.15) gives an estimate of the number of steps they take to arrive, within this model.

Equation (18.15) is rather similar to the corresponding expression for the model of Kleinberg, Eq. (18.8), which is not a coincidence since the mechanisms by which the message-passing process proceeds are similar in the two cases. Taking the limit of large n and making use of Eq. (18.14), we find that

$$\ell = \begin{cases} D(n/g)^{1-\beta} & \text{for } \beta < 1, \\ E \log^2(n/g) & \text{for } \beta = 1, \\ F(n/g)^{\beta-1} & \text{for } \beta > 1, \end{cases} \quad (18.16)$$

where D , E , and F are constants.

These results have the same functional form as those of Eq. (18.9) for Kleinberg's model and tell us that it is indeed possible for Milgram's experiment to succeed in networks of this type, but only for the special parameter value $\beta = 1$. For all other values, the number of steps ℓ taken to reach the target increases as a power of n . Note that $\beta = 1$ is precisely the point at which the expected number of acquaintances is the same at all distances, since this is the point at which Eq. (18.11) becomes independent of m . In other words, the Milgram experiment succeeds only when the decrease with distance in the probability of knowing a particular person is exactly canceled out by the increase in the number of people there are to know.

The model of Watts *et al.* confirms Kleinberg's results in the context of a more realistic network. The results are, however, somewhat mysterious in a way. The idea that the network must be tuned to a special point in order for the Milgram experiment to succeed is surprising. The Milgram experiment does appear to succeed when conducted on real-world social networks, but on the face of it there is no clear reason why real-world networks should fall at this special point. Is it really true that if the world happened to be a little different from the way it is, the experiment would fail? This is a point that is not yet fully understood. It is possible that the models are missing some important feature of real-world network structure that makes message passing more robust and less dependent on the precise tuning of the network. Or perhaps people are

using a different scheme for passing messages that works substantially better than our greedy algorithm. We have assumed, for instance, that if a person does not know someone closer to the target than themselves then they pass the message at random to someone in their own subtree. However, it's reasonable to suppose that real people might use a smarter algorithm. If you have a message destined for someone in a particular country, say, but you don't know anyone in that country, then you might strategically pass the message to an acquaintance who you know does have contacts there. Such a tactic could significantly reduce the number of steps the message takes.

It is also possible, however, that our model is basically correct but that the world is in fact only rather loosely tuned to the special point $\beta = 1$. For values of β close to 1 the power of n in Eq. (18.15) is small and hence ℓ still grows quite slowly. Indeed it is in general difficult to distinguish experimentally between low powers and logarithms, so any value of β in the rough vicinity of $\beta = 1$ could result in good apparent performance in the message-passing experiment.

EXERCISES

18.1 Suppose that we use a web crawler to crawl a small portion of the Web, starting from a randomly chosen web page somewhere in the large in-component. Let us model the crawl as a breadth-first search starting from the given node and proceeding for r "waves" of search, i.e., until it reaches nodes that are r steps away from the start. Let S_i be the size of the large in-component.

- a) What is the probability that a given web page has been crawled at the "zeroth" wave of the algorithm, i.e., when only the one starting page has been crawled?
- b) Argue that the probability p_i that a page is first reached on the r th wave is given approximately by $\mathbf{p}(r) = \mathbf{A}\mathbf{p}(r - 1)$, where $\mathbf{p} = (p_1, p_2, \dots)$. Why is this relation only approximate in general?
- c) Hence argue that the probability of a page being found in a small crawl is roughly proportional to the eigenvector centrality of the page. Recall that the eigenvector centrality is zero for nodes in the in-component that don't also belong to the strongly connected component (see Section 7.1.2). Explain why this makes sense in the present context.

18.2 Suppose that a search is performed on a peer-to-peer network using the following algorithm. Each node on the network maintains a record of the items held by each of its neighbors. The node originating a search queries one of its neighbors, chosen uniformly at random, for a desired item and the neighbor responds either that it or one

of its neighbors has the item, in which case the search ends, or that they do not. In the latter case, the neighboring node then passes the query on to one of *its* neighbors, also chosen at random, and the process repeats until the item is found. Effectively, therefore, the search query makes a random walk on the network.

- a) Argue that, in the limit of a large number of steps, the probability that the query encounters a node i on any particular step is $k_i/2m$, where k_i is the degree as usual and m is the total number of edges in the network.
- b) Upon arriving at a node of degree k , the search learns (at most) about the items held by all of that node's k neighbors except for the one the query is coming from, for a total of $k - 1$ nodes. Show that on average at each step the search learns about the contents of approximately $\langle k^2 \rangle / \langle k \rangle - 1$ nodes and hence that, for a target item that can be found at a fraction c of the nodes in the network, the expected number of copies of the item found on a given step is $c(\langle k^2 \rangle / \langle k \rangle - 1)$.
- c) Argue that the probability of not finding the target item on any particular step is approximately $q = \exp[c(1 - \langle k^2 \rangle / \langle k \rangle)]$ and that the average number of steps it takes to find a copy of the item is $1/(1 - q)$.
- d) On a network with a power-law degree distribution with exponent less than 3, so that $\langle k^2 \rangle \rightarrow \infty$, this last result implies that in the limit of large network size the search should end after only one step. Is this really true? If not, explain why not.

Although the random walk is not a realistic model for actual network search it is nonetheless useful: presumably more intelligent search strategies will find results more quickly than a mindless random walk and hence the random walk provides an upper bound on the length of the search needed to find an item. In particular, if the random walk works well, as in the example above, then it suggests that more intelligent forms of search will also work well.

18.3 The network navigation model of Kleinberg described in Section 18.3.1 is a one-dimensional version of what was, originally, a two-dimensional model. In Kleinberg's original version, the model was built on a two-dimensional square lattice with nodes connected by shortcuts with probability proportional to $r^{-\alpha}$ where r is the "Manhattan" distance between the nodes, i.e., the network distance in terms of number of edges traversed (rather than the Euclidean distance). Following the outline of Section 18.3.1, sketch an argument to show for this variant of the model that it is possible to find a target node in $O(\log^2 n)$ steps, but only if $\alpha = 2$.

REFERENCES

- [1] Abello, J., Buchsbaum, A., and Westbrook, J., A functional approach to external graph algorithms, in *Proceedings of the 6th European Symposium on Algorithms*, Springer, Berlin (1998).
- [2] Abramowitz, M. and Stegun, I. A., eds., *Handbook of Mathematical Functions*, Dover Publishing, New York (1974).
- [3] Achlioptas, D., Clauset, A., Kempe, D., and Moore, C., On the bias of traceroute sampling, in *Proceedings of the 37th ACM Symposium on Theory of Computing*, Association of Computing Machinery, New York (2005).
- [4] Adamic, L. A. and Glance, N., The political blogosphere and the 2004 US election, in *Proceedings of the WWW-2005 Workshop on the Weblogging Ecosystem*, Association of Computing Machinery, New York (2005).
- [5] Adamic, L. A. and Huberman, B. A., The nature of markets in the World Wide Web, *Quarterly Journal of Electronic Commerce* **1**, 512 (2000).
- [6] Adamic, L. A., Lukose, R. M., Puniyani, A. R., and Huberman, B. A., Search in power-law networks, *Phys. Rev. E* **64**, 046135 (2001).
- [7] Adler, J., Bootstrap percolation, *Physica A* **171**, 453–470 (1991).
- [8] Ahn, Y.-Y., Bagrow, J. P., and Lehmann, S., Link communities reveal multiscale complexity in networks, *Nature* **466**, 761–764 (2010).
- [9] Ahuja, R. K., Magnanti, T. L., and Orlin, J. B., *Network Flows: Theory, Algorithms, and Applications*, Prentice Hall, Upper Saddle River, NJ (1993).
- [10] Aiello, W., Chung, F., and Lu, L., A random graph model for massive graphs, in *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, pp. 171–180, Association of Computing Machinery, New York (2000).
- [11] Aiello, W., Chung, F., and Lu, L., Random evolution of massive graphs, in J. Abello, P. M. Pardalos, and M. G. C. Resende, eds., *Handbook of Massive Data Sets*, pp. 97–122, Kluwer, Dordrecht (2002).
- [12] Airoldi, E. M., Blei, D. M., Fienberg, S. E., and Xing, E. P., Mixed membership stochastic blockmodels, *J. Mach. Learn. Res.* **9**, 1981–2014 (2008).
- [13] Albert, R., Albert, I., and Nakarado, G. L., Structural vulnerability of the North American power grid, *Phys. Rev. E* **69**, 025103 (2004).
- [14] Albert, R. and Barabási, A.-L., Topology of evolving networks: Local events and universality, *Phys. Rev. Lett.* **85**, 5234–5237 (2000).
- [15] Albert, R. and Barabási, A.-L., Statistical mechanics of complex networks, *Rev. Mod. Phys.* **74**, 47–97 (2002).
- [16] Albert, R., Jeong, H., and Barabási, A.-L., Diameter of the world-wide web, *Nature* **401**, 130–131 (1999).
- [17] Albert, R., Jeong, H., and Barabási, A.-L., Attack and error tolerance of complex networks, *Nature* **406**, 378–382 (2000).
- [18] Aldous, D. J., Spatial transportation networks with transfer costs: Asymptotic optimality of hub-and-spoke models, *Math. Proc. Camb. Phil. Soc.* **145**, 471–487 (2008).
- [19] Ali, I., Cook, W. D., and Kress, M., On the minimum violations ranking of a tournament, *Manag. Sci.* **32**, 660–672 (1986).
- [20] Amaral, L. A. N., Scala, A., Barthélemy, M., and Stanley, H. E., Classes of small-world networks, *Proc. Natl. Acad. Sci. USA* **97**, 11149–11152 (2000).
- [21] Anderson, R. M. and May, R. M., *Infectious Diseases of Humans*, Oxford University Press, Oxford (1991).
- [22] Anderson, W. N. and Morley, T. D., Eigenvalues of the Laplacian of a graph, *Linear and Multilinear Algebra* **18**, 141–145 (1985).

- [23] Anthonisse, J. M., The rush in a directed graph, Technical Report BN 9/71, Stichting Mathematisch Centrum, Amsterdam (1971).
- [24] Appel, K. and Haken, W., Every planar map is four colorable. II: Reducibility, *Illinois J. Math.* **21**, 491–567 (1977).
- [25] Appel, K. and Haken, W., The solution of the four-color map problem, *Sci. Am.* **237**, 108–121 (1977).
- [26] Appel, K., Haken, W., and Koch, J., Every planar map is four colorable. I: Discharging, *Illinois J. Math.* **21**, 429–490 (1977).
- [27] Appleby, M. C., Social rank and food access in red deer stags, *Behaviour* **74**, 294–309 (1980).
- [28] Aral, S. and Walker, D., Identifying influential and susceptible members of social networks, *Science* **337**, 337–341 (2012).
- [29] Aref, S. and Wilson, M. C., Measuring partial balance in signed networks, preprint arxiv:1509.04037 (2015).
- [30] Arenas, A., Díaz-Guilera, A., Kurths, J., Moreno, Y., and Zhou, C., Synchronization in complex networks, *Phys. Rep.* **469**, 93–153 (2008).
- [31] Arianos, S., Bompard, E., Carbone, A., and Xue, F., Power grid vulnerability: A complex network approach, *Chaos* **19**, 013119 (2009).
- [32] Auerbach, F., Das Gesetz der Bevölkerungskonzentration, *Petermanns Geographische Mitteilungen* **59**, 74–76 (1913).
- [33] Axelrod, R. and Bennett, D. S., A landscape theory of aggregation, *Br. J. Polit. Sci.* **23**, 211–233 (1993).
- [34] Bagler, G., Analysis of the airport network of India as a complex weighted network, *Physica A* **387**, 2972–2980 (2008).
- [35] Bagler, G. and Sinha, S., Network properties of protein structures, *Physica A* **346**, 27–33 (2005).
- [36] Bailey, N. T. J., *The Mathematical Theory of Infectious Diseases and Its Applications*, Hafner Press, New York (1975).
- [37] Ball, B., Karrer, B., and Newman, M. E. J., An efficient and principled method for detecting communities in networks, *Phys. Rev. E* **84**, 036103 (2011).
- [38] Ball, B. and Newman, M. E. J., Friendship networks and social status, *Netw. Sci.* **1**, 16–30 (2013).
- [39] Banavar, J. R., Maritan, A., and Rinaldo, A., Size and form in efficient transportation networks, *Nature* **399**, 130–132 (1999).
- [40] Barabási, A.-L. and Albert, R., Emergence of scaling in random networks, *Science* **286**, 509–512 (1999).
- [41] Barabási, A.-L., Albert, R., and Jeong, H., Scale-free characteristics of random networks: The topology of the World Wide Web, *Physica A* **281**, 69–77 (2000).
- [42] Barabási, A.-L., Gulbahce, N., and Loscalzo, J., Network medicine: A network-based approach to human disease, *Nat. Rev. Genet.* **12**, 57–68 (2011).
- [43] Barabási, A.-L., Jeong, H., Ravasz, E., Nédá, Z., Schuberts, A., and Vicsek, T., Evolution of the social network of scientific collaborations, *Physica A* **311**, 590–614 (2002).
- [44] Barbour, A. and Mollison, D., Epidemics and random graphs, in J. P. Gabriel, C. Lefevre, and P. Picard, eds., *Stochastic Processes in Epidemic Theory*, pp. 86–89, Springer, New York (1990).
- [45] Barrat, A., Barthélemy, M., and Vespignani, A., *Dynamical Processes on Complex Networks*, Cambridge University Press, Cambridge (2008).
- [46] Barthélemy, M., Spatial networks, *Phys. Rep.* **499**, 1–101 (2011).
- [47] Barthélemy, M., Barrat, A., Pastor-Satorras, R., and Vespignani, A., Velocity and hierarchical spread of epidemic outbreaks in scale-free networks, *Phys. Rev. Lett.* **92**, 178701 (2004).
- [48] Barthélemy, M., Barrat, A., Pastor-Satorras, R., and Vespignani, A., Dynamical patterns of epidemic outbreaks in complex heterogeneous networks, *J. Theor. Bio.* **235**, 275–288 (2005).
- [49] Batson, J., Spielman, D. A., Srivastava, N., and Teng, S.-H., Spectral sparsification of graphs: Theory and algorithms, *Commun. ACM* **56**, 87–94 (2013).
- [50] Baxter, G. J., Dorogovtsev, S. N., Goltsev, A. V., and Mendes, J. F. F., Bootstrap percolation on complex networks, *Phys. Rev. E* **82**, 011103 (2010).
- [51] Bearman, P., Moody, J., and Faris, R., Networks and history, *Complexity* **8**, 61–71 (2003).
- [52] Bearman, P. S., Moody, J., and Stovel, K., Chains of affection: The structure of adolescent romantic and sexual networks, *Am. J. Sociol.* **110**, 44–91 (2004).
- [53] Berger, S. I. and Iyengar, R., Network analyses in systems pharmacology, *Bioinformatics* **25**, 2466–2472 (2009).
- [54] Bernard, H. R., Johnsen, E. C., Killworth, P. D., and Robinson, S., Estimating the size of an average

REFERENCES

- personal network and of an event population, in M. Kochen, ed., *The Small World*, pp. 159–175, Ablex Publishing, Norwood, NJ (1989).
- [55] Bernard, H. R., Johnsen, E. C., Killworth, P. D., and Robinson, S., Estimating the size of an average personal network and of an event population: Some empirical results, *Social Sci. Res.* **20**, 109–121 (1991).
- [56] Bernard, H. R. and Killworth, P. D., Informant accuracy in social network data II, *Human Communications Res.* **4**, 3–18 (1977).
- [57] Bernard, H. R., Killworth, P. D., Evans, M. J., McCarty, C., and Shelley, G. A., Studying social relations cross-culturally, *Ethnology* **2**, 155–179 (1988).
- [58] Bernard, H. R., Killworth, P. D., and Sailer, L., Informant accuracy in social network data IV: A comparison of clique-level structure in behavioral and cognitive network data, *Soc. Networks* **2**, 191–218 (1980).
- [59] Bernard, H. R., Killworth, P. D., and Sailer, L., Informant accuracy in social network data V: An experimental attempt to predict actual communication from recall data, *Social Sci. Res.* **11**, 30–66 (1982).
- [60] Bianconi, G., *Multilayer Networks: Structure and Function*, Oxford University Press, Oxford (2018).
- [61] Bianconi, G. and Capocci, A., Number of loops of size h in growing scale-free networks, *Phys. Rev. Lett.* **90**, 078701 (2003).
- [62] Bickel, P. J. and Chen, A., A nonparametric view of network models and Newman–Girvan and other modularities, *Proc. Natl. Acad. Sci. USA* **106**, 21068–21073 (2009).
- [63] Blei, D. M., Ng, A. Y., and Jordan, M. I., Latent Dirichlet allocation, *J. Mach. Learn. Res.* **3**, 993–1022 (2003).
- [64] Blondel, V. D., Decuyper, A., and Krings, G., A survey of results of mobile phone datasets analysis, *EPJ Data Sci.* **4**, 10 (2015).
- [65] Blondel, V. D., Gajardo, A., Heymans, M., Senellart, P., and Dooren, P. V., A measure of similarity between graph vertices: Applications to synonym extraction and web searching, *SIAM Rev.* **46**, 647–666 (2004).
- [66] Blondel, V. D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E., Fast unfolding of communities in large networks, *J. Stat. Mech.* **2008**, P10008 (2008).
- [67] Boccaletti, S., Bianconi, G., Criado, R., del Genio, C. I., Gomez-Gardenes, J., Romance, M., Sendina-Nadal, I., Wang, Z., and Zanin, M., The structure and dynamics of multilayer networks, *Phys. Rep.* **544**, 1–122 (2014).
- [68] Boccaletti, S., Latora, V., Moreno, Y., Chavez, M., and Hwang, D.-U., Complex networks: Structure and dynamics, *Phys. Rep.* **424**, 175–308 (2006).
- [69] Bollobás, B., A probabilistic proof of an asymptotic formula for the number of labelled regular graphs, *Eur. J. Combinatorics* **1**, 311–316 (1980).
- [70] Bollobás, B., *Random Graphs*, 2nd edn., Academic Press, New York (2001).
- [71] Bollobás, B. and Riordan, O., Sparse graphs: Metrics and random models, *Random Struct. Alg.* **39**, 1–38 (2010).
- [72] Bollobás, B., Riordan, O., Spencer, J., and Tusnády, G., The degree sequence of a scale-free random graph process, *Random Struct. Alg.* **18**, 279–290 (2001).
- [73] Bonacich, P. F., Power and centrality: A family of measures, *Am. J. Sociol.* **92**, 1170–1182 (1987).
- [74] Bond, R. M., Fariss, C. J., Jones, J. J., Kramer, A. D. I., Marlow, C., Settle, J. E., and Fowler, J. H., A 61-million-person experiment in social influence and political mobilization, *Nature* **489**, 295–298 (2012).
- [75] Borgatti, S. P., Structural holes: Unpacking Burt’s redundancy measures, *Connections* **20**(1), 35–38 (1997).
- [76] Borgatti, S. P., Centrality and network flow, *Soc. Networks* **27**, 55–71 (2005).
- [77] Borgatti, S. P., Carley, K. M., and Krackhardt, D., On the robustness of centrality measures under conditions of imperfect data, *Soc. Networks* **28**, 124–136 (2006).
- [78] Borgatti, S. P. and Everett, M. G., Models of core/periphery structures, *Soc. Networks* **21**, 375–395 (1999).
- [79] Borgatti, S. P., Mehra, A., Brass, D. J., and Labianca, G., Network analysis in the social sciences, *Science* **323**, 892–895 (2009).
- [80] Borodin, A., Roberts, G. O., Rosenthal, J. S., and Tsaparas, P., Finding authorities and hubs from link structures on the World Wide Web, in V. Y. Shen, N. Saito, M. R. Lyu, and M. E. Zurko, eds., *Proceedings of the 10th International World Wide Web*

- Conference, pp. 415–429, Association of Computing Machinery, New York (2001).
- [81] Brandes, U., Delling, D., Gaertler, M., Görke, R., Hofer, M., Nikoloski, Z., and Wagner, D., On finding graph clusterings with maximum modularity, in *Proceedings of the 33rd International Workshop on Graph-Theoretic Concepts in Computer Science*, no. 4769 in Lecture Notes in Computer Science, Springer, Berlin (2007).
- [82] Brin, S. and Page, L., The anatomy of a large-scale hypertextual Web search engine, *Comput. Netw.* **30**, 107–117 (1998).
- [83] Brinda, K. V. and Vishveshwara, S., A network representation of protein structures: Implications for protein stability, *Biophys. J.* **89**, 4159–4170 (2005).
- [84] Broder, A., Kumar, R., Maghoul, F., Raghavan, P., Rajagopalan, S., Stata, R., Tomkins, A., and Wiener, J., Graph structure in the web, *Comput. Netw.* **33**, 309–320 (2000).
- [85] Broido, A. and Claffy, K. C., Internet topology: Connectivity of IP graphs, in S. Fahmy and K. Park, eds., *Scalability and Traffic Control in IP Networks*, no. 4526 in Proc. SPIE, pp. 172–187, International Society for Optical Engineering, Bellingham, WA (2001).
- [86] Bullmore, E. T. and Bassett, D. S., Brain graphs: Graphical models of the human brain connectome, *Annu. Rev. Clin. Psychol.* **7**, 113–140 (2011).
- [87] Burlando, B., The fractal dimension of taxonomic systems, *J. Theor. Bio.* **146**, 99–114 (1990).
- [88] Burt, R. S., Network items and the General Social Survey, *Soc. Networks* **6**, 293–339 (1984).
- [89] Burt, R. S., *Structural Holes: The Social Structure of Competition*, Harvard University Press, Cambridge, MA (1992).
- [90] Butts, C. T., Network inference, error, and informant (in)accuracy: A Bayesian approach, *Soc. Networks* **25**, 103–140 (2003).
- [91] Calabrese, F., Smoreda, Z., Blondel, V. D., and Ratti, C., Interplay between telecommunications and face-to-face interactions: A study using mobile phone data, *PLOS One* **6**, e20814 (2011).
- [92] Caldarelli, G., Pastor-Satorras, R., and Vespignani, A., Structure of cycles and local ordering in complex networks, *Eur. Phys. J. B* **38**, 183–186 (2004).
- [93] Callaway, D. S., Newman, M. E. J., Strogatz, S. H., and Watts, D. J., Network robustness and fragility: Percolation on random graphs, *Phys. Rev. Lett.* **85**, 5468–5471 (2000).
- [94] Cano, P., Celma, O., Koppenberger, M., and Buldú, J. M., Topology of music recommendation networks, *Chaos* **16**, 013107 (2006).
- [95] Cardillo, A., Scellato, S., Latora, V., and Porta, S., Structural properties of planar graphs of urban street patterns, *Phys. Rev. E* **73**, 066107 (2006).
- [96] Carvalho, R., Buzna, L., Bono, F., Gutierrez, E., Just, W., and Arrowsmith, D., Robustness of trans-European gas networks, *Phys. Rev. E* **80**, 016106 (2009).
- [97] Catania, J. A., Coates, T. J., Kegels, S., and Fullilove, M. T., The population-based AMEN (AIDS in Multi-Ethnic Neighborhoods) study, *Am. J. Public Health* **82**, 284–287 (1992).
- [98] Centola, D. and Macy, M., Complex contagions and the weakness of long ties, *Am. J. Sociology* **113**, 702–734 (2007).
- [99] Chalupa, J., Leath, P. L., and Reich, G. R., Bootstrap percolation on a Bethe lattice, *J. Phys. C* **12**, L31–35 (1979).
- [100] Chapelle, O., Schölkopf, B., and Zien, A., eds., *Semi-Supervised Learning*, MIT Press, Cambridge, MA (2006).
- [101] Chen, P. and Redner, S., Community structure of the Physical Review citation network, *J. Informetr.* **4**, 278–290 (2010).
- [102] Chen, Q., Chang, H., Govindan, R., Jamin, S., Shenker, S. J., and Willinger, W., The origin of power laws in Internet topologies revisited, in *Proceedings of the 21st Annual Joint Conference of the IEEE Computer and Communications Societies*, IEEE Computer Society, New York (2002).
- [103] Chung, F. and Lu, L., The average distances in random graphs with given expected degrees, *Proc. Natl. Acad. Sci. USA* **99**, 15879–15882 (2002).
- [104] Chung, F., Lu, L., and Vu, V., Spectra of random graphs with given expected degrees, *Proc. Natl. Acad. Sci. USA* **100**, 6313–6318 (2003).
- [105] Chung, K. and Deisseroth, K., CLARITY for mapping the nervous system, *Nat. Methods* **10**, 508–513 (2013).

REFERENCES

- [106] Clarkson, G. and DeKorte, D., The problem of patent thickets in convergent technologies, in W. S. Bainbridge and M. C. Roco, eds., *Progress in Convergence: Technologies for Human Wellbeing*, no. 1093 in *Annals of the New York Academy of Science*, pp. 180–200, New York Academy of Sciences, New York (2006).
- [107] Clauset, A., Arbesman, S., and Larremore, D. B., Systematic inequality and hierarchy in faculty hiring networks, *Sci. Adv.* **1**, e1400005 (2015).
- [108] Clauset, A. and Moore, C., Accuracy and scaling phenomena in Internet mapping, *Phys. Rev. Lett.* **94**, 018701 (2005).
- [109] Clauset, A., Moore, C., and Newman, M. E. J., Hierarchical structure and the prediction of missing links in networks, *Nature* **453**, 98–101 (2008).
- [110] Clauset, A., Newman, M. E. J., and Moore, C., Finding community structure in very large networks, *Phys. Rev. E* **70**, 066111 (2004).
- [111] Clauset, A., Shalizi, C. R., and Newman, M. E. J., Power-law distributions in empirical data, *SIAM Rev.* **51**, 661–703 (2009).
- [112] Cohen, J. E., *Ecologists' Co-operative Web Bank, Version 1.0: Machine-Readable Data Base of Food Webs*, Rockefeller University, New York (1989).
- [113] Cohen, R., Erez, K., ben-Avraham, D., and Havlin, S., Resilience of the Internet to random breakdowns, *Phys. Rev. Lett.* **85**, 4626–4628 (2000).
- [114] Cohen, R. and Havlin, S., Scale-free networks are ultrasmall, *Phys. Rev. Lett.* **90**, 058701 (2003).
- [115] Cohen, R., Havlin, S., and ben-Avraham, D., Efficient immunization strategies for computer networks and populations, *Phys. Rev. Lett.* **91**, 247901 (2003).
- [116] Cole, B. J., Dominance hierarchies in Leptothorax ants, *Science* **212**, 83–84 (1981).
- [117] Coleman, J. S., Katz, E., and Menzel, H., The diffusion of an innovation among physicians, *Sociometry* **20**, 253–270 (1957).
- [118] Colizza, V., Barrat, A., Barthélemy, M., and Vespignani, A., The role of the airline transportation network in the prediction and predictability of global epidemics, *Proc. Natl. Acad. Sci. USA* **103**, 2015–2020 (2006).
- [119] Côme, E. and Latouche, P., Model selection and clustering in stochastic block models based on the exact integrated complete data likelihood, *Stat. Model.* **15**, 564–589 (2015).
- [120] Condon, A. and Karp, R. M., Algorithms for graph partitioning on the planted partition model, *Random Struct. Alg.* **18**, 116–140 (2001).
- [121] Connor, R. C., Heithaus, M. R., and Barre, L. M., Superalliance of bottlenose dolphins, *Nature* **397**, 571–572 (1999).
- [122] Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C., *Introduction to Algorithms*, 2nd edn., MIT Press, Cambridge, MA (2001).
- [123] Cox, R. A. K., Felton, J. M., and Chung, K. C., The concentration of commercial success in popular music: An analysis of the distribution of gold records, *J. Cult. Econ.* **19**, 333–340 (1995).
- [124] Crovella, M. E. and Bestavros, A., Self-similarity in World Wide Web traffic: Evidence and possible causes, in B. E. Gaither and D. A. Reed, eds., *Proceedings of the 1996 ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, pp. 148–159, Association of Computing Machinery, New York (1996).
- [125] Crucitti, P., Latora, V., and Marchiori, M., A topological analysis of the Italian electric power grid, *Physica A* **338**, 92–97 (2004).
- [126] Csermely, P., London, A., Wu, L.-Y., and Uzzi, B., Structure and dynamics of core/periphery networks, *J. Complex Netw.* **1**, 93–123 (2013).
- [127] D'Angelo, C. A., Giuffrida, C., and Abramo, G., A heuristic approach to author name disambiguation in bibliometrics databases for large-scale research assessments, *J. Assoc. Inf. Sci. Technol.* **62**, 257–269 (2011).
- [128] Danon, L., Duch, J., Diaz-Guilera, A., and Arenas, A., Comparing community structure identification, *J. Stat. Mech.* **2005**, P09008 (2005).
- [129] Davis, A., Gardner, B. B., and Gardner, M. R., *Deep South*, University of Chicago Press, Chicago (1941).
- [130] Davis, G. F. and Greve, H. R., Corporate elite networks and governance changes in the 1980s, *Am. J. Sociol.* **103**, 1–37 (1997).
- [131] Davis, G. F., Yoo, M., and Baker, W. E., The small world of the American corporate elite, 1982–2001, *Strateg. Organ.* **1**, 301–326 (2003).
- [132] Davis, J. A., Clustering and structural balance in graphs, *Human Relations* **20**, 181–187 (1967).

- [133] de Castro, R. and Grossman, J. W., Famous trails to Paul Erdős, *Math. Intelligencer* **21**, 51–63 (1999).
- [134] De Domenico, M., Granell, C., Porter, M. A., and Arenas, A., The physics of multilayer networks, *Nat. Phys.* **12**, 901–906 (2016).
- [135] De Domenico, M., Solé-Ribalta, A., Gómez, S., and Arenas, A., Navigability of interconnected networks under random failures, *Proc. Natl. Acad. Sci. USA* **111**, 8351–8356 (2014).
- [136] De Vries, H., Finding a dominance order most consistent with a linear hierarchy: A new procedure and review, *Animal Behav.* **55**, 827–843 (1998).
- [137] De Vries, H., Stevens, J. M. G., and Vervaecke, H., Measuring and testing the steepness of dominance hierarchies, *Animal Behav.* **55**, 585–592 (2006).
- [138] Decelle, A., Krzakala, F., Moore, C., and Zdeborová, L., Asymptotic analysis of the stochastic block model for modular networks and its algorithmic applications, *Phys. Rev. E* **84**, 066106 (2011).
- [139] Decelle, A., Krzakala, F., Moore, C., and Zdeborová, L., Inference and phase transitions in the detection of modules in sparse networks, *Phys. Rev. Lett.* **107**, 065701 (2011).
- [140] Dobson, I., Carreras, B. A., Lynch, V. E., and Newman, D. E., Complex systems analysis of series of blackouts: Cascading failure, critical points, and self-organization, *Chaos* **17**, 026103 (2007).
- [141] Dodds, P. S., Harris, K. D., Kloumann, I. M., Bliss, C. A., and Danforth, C. M., Temporal patterns of happiness and information in a global social network: Hedonometrics and Twitter, *PLOS One* **6**, e26752 (2011).
- [142] Dodds, P. S., Muhamad, R., and Watts, D. J., An experimental study of search in global social networks, *Science* **301**, 827–829 (2003).
- [143] Dodds, P. S. and Rothman, D. H., Geometry of river networks, *Phys. Rev. E* **63**, 016115, 016116, & 016117 (2001).
- [144] Dorogovtsev, S. N., Goltsev, A. V., and Mendes, J. F. F., Ising model on networks with an arbitrary distribution of connections, *Phys. Rev. E* **66**, 016104 (2002).
- [145] Dorogovtsev, S. N. and Mendes, J. F. F., Scaling behaviour of developing and decaying networks, *Europhys. Lett.* **52**, 33–39 (2000).
- [146] Dorogovtsev, S. N. and Mendes, J. F. F., Language as an evolving word web, *Proc. R. Soc. London B* **268**, 2603–2606 (2001).
- [147] Dorogovtsev, S. N. and Mendes, J. F. F., Evolution of networks, *Adv. Phys.* **51**, 1079–1187 (2002).
- [148] Dorogovtsev, S. N., Mendes, J. F. F., and Samukhin, A. N., Structure of growing networks with preferential linking, *Phys. Rev. Lett.* **85**, 4633–4636 (2000).
- [149] Dorogovtsev, S. N., Mendes, J. F. F., and Samukhin, A. N., Giant strongly connected component of directed networks, *Phys. Rev. E* **64**, 025101 (2001).
- [150] Drews, C., The concept and definition of dominance in animal behaviour, *Behaviour* **125**, 283–313 (1993).
- [151] Duch, J. and Arenas, A., Community detection in complex networks using extremal optimization, *Phys. Rev. E* **72**, 027104 (2005).
- [152] Dunne, J. A., Labandeira, C. C., and Williams, R. J., Highly resolved early Eocene food webs show development of modern trophic structure after the end-Cretaceous extinction, *Proc. R. Soc. London B* **281**, 20133280 (2014).
- [153] Dunne, J. A., Williams, R. J., and Martinez, N. D., Food-web structure and network theory: The role of connectance and size, *Proc. Natl. Acad. Sci. USA* **99**, 12917–12922 (2002).
- [154] Eagle, N. and Pentland, A., Reality mining: Sensing complex social systems, *J. Pers. Ubiquitous Comput.* **10**, 255–268 (2006).
- [155] Eagle, N., Pentland, A., and Lazer, D., Inferring friendship network structure by using mobile phone data, *Proc. Natl. Acad. Sci. USA* **106**, 15274–15278 (2009).
- [156] Ebel, H., Mielsch, L.-I., and Bornholdt, S., Scale-free topology of e-mail networks, *Phys. Rev. E* **66**, 035103 (2002).
- [157] Eckmann, J.-P. and Moses, E., Curvature of co-links uncovers hidden thematic layers in the World Wide Web, *Proc. Natl. Acad. Sci. USA* **99**, 5825–5829 (2002).
- [158] Erdős, P. and Rényi, A., On random graphs, *Publ. Math.* **6**, 290–297 (1959).
- [159] Erdős, P. and Rényi, A., On the evolution of random graphs, *Publications of the Mathematical Institute of the Hungarian Academy of Sciences* **5**, 17–61 (1960).

REFERENCES

- [160] Erdős, P. and Rényi, A., On the strength of connect- edness of a random graph, *Acta Math. Sci. Hungary* **12**, 261–267 (1961).
- [161] Érdi, P., Makovi, K., Somogyvári, Z., Strandburg, K., Tobochnik, J., Volf, P., and Zláányi, L., Predic- tion of emerging technologies based on analysis of the US patent citation network, *Scientometrics* **95**, 225–242 (2013).
- [162] Erickson, B., Some problems of inference from chain data, in K. F. Schuessler, ed., *Sociological Methodology 1979*, pp. 276–302, Jossey-Bass, San Francisco (1978).
- [163] Erman, N. and Todorovski, L., The effects of mea- surement error in case of scientific network analy- sis, *Scientometrics* **104**, 453–473 (2015).
- [164] Estoup, J. B., *Gammes Stenographiques*, Institut Stenographique de France, Paris (1916).
- [165] Eubank, S., Guclu, H., Kumar, V. S. A., Marathe, M. V., Srinivasan, A., Toroczkai, Z., and Wang, N., Modelling disease outbreaks in realistic urban social networks, *Nature* **429**, 180–184 (2004).
- [166] Evans, T. S. and Lambiotte, R., Line graphs, link partitions, and overlapping communities, *Phys. Rev. E* **80**, 016105 (2009).
- [167] Facchetti, G., Iacono, G., and Altafini, C., Comput- ing global structural balance in large-scale signed social networks, *Proc. Natl. Acad. Sci. USA* **108**, 20953–20958 (2011).
- [168] Faloutsos, M., Faloutsos, P., and Faloutsos, C., On power-law relationships of the internet topology, *Comput. Commun. Rev.* **29**, 251–262 (1999).
- [169] Fararo, T. J. and Sunshine, M., *A Study of a Biased Friendship Network*, Syracuse University Press, Syra- cuse (1964).
- [170] Feld, S., Why your friends have more friends than you do, *Am. J. Sociol.* **96**, 1464–1477 (1991).
- [171] Feld, S. L. and Carter, W. C., Detecting mea- surement bias in respondent reports of personal net- works, *Soc. Networks* **24**, 365–383 (2002).
- [172] Fernholz, D. and Ramachandran, V., The diameter of sparse random graphs, *Random Struct. Alg.* **31**, 482–516 (2007).
- [173] Ferreira, A. A., Goncalves, M. A., and Laender, A. H. F., A brief survey of automatic methods for au- thor name disambiguation, *SIGMOD Record* **41**, 15–26 (2012).
- [174] Ferrer i Cancho, R., Janssen, C., and Solé, R. V., Topology of technology graphs: Small world pat- terns in electronic circuits, *Phys. Rev. E* **64**, 046119 (2001).
- [175] Ferrer i Cancho, R. and Solé, R. V., The small world of human language, *Proc. R. Soc. London B* **268**, 2261–2265 (2001).
- [176] Ferrer i Cancho, R. and Solé, R. V., Optimiza- tion in complex networks, in R. Pastor-Satorras, J. Rubi, and A. Díaz-Guilera, eds., *Statistical Mech- anics of Complex Networks*, no. 625 in Lecture Notes in Physics, pp. 114–125, Springer, Berlin (2003).
- [177] Fiedler, M., Algebraic connectivity of graphs, *Czech. Math. J.* **23**, 298–305 (1973).
- [178] Fields, S. and Song, O., A novel genetic system to detect protein-protein interactions, *Nature* **340**, 245–246 (1989).
- [179] Fisher, M. E. and Essam, J. W., Some cluster size and percolation problems, *J. Math. Phys.* **2**, 609–619 (1961).
- [180] Flack, J. C., Girvan, M., de Waal, F. B. M., and Krakauer, D. C., Policing stabilizes construction of social niches in primates, *Nature* **439**, 426–429 (2006).
- [181] Flake, G. W., Lawrence, S. R., Giles, C. L., and Co- etzee, F. M., Self-organization and identification of Web communities, *IEEE Computer* **35**, 66–71 (2002).
- [182] Flory, P. J., Molecular size distribution in three di- mensional polymers. I: Gelation, *J. Am. Chem. Soc.* **63**, 3083–3090 (1941).
- [183] Fortunato, S., Community detection in graphs, *Phys. Rep.* **486**, 75–174 (2010).
- [184] Fortunato, S. and Barthélemy, M., Resolution limit in community detection, *Proc. Natl. Acad. Sci. USA* **104**, 36–41 (2007).
- [185] Fortunato, S. and Hric, D., Community detection in networks: A user guide, *Phys. Rep.* **659**, 1–44 (2016).
- [186] Fowler, J. H. and Jeon, S., The authority of Supreme Court precedent, *Soc. Networks* **30**, 16–30 (2008).
- [187] Fowler, J. H., Johnson, T. R., Spriggs II, J. F., Jeon, S., and Wahlbeck, P. J., Network analysis and the law: Measuring the legal importance of Supreme Court precedents, *Political Anal.* **15**, 324–346 (2007).
- [188] Frank, O., Estimation of population totals by use of snowball samples, in P. W. Holland and S. Lein-

- hardt, eds., *Perspectives on Social Network Research*, pp. 319–348, Academic Press, New York (1979).
- [189] Freeman, L. C., A set of measures of centrality based upon betweenness, *Sociometry* **40**, 35–41 (1977).
- [190] Freeman, L. C., Finding social groups: A meta-analysis of the southern women data, in R. Breiger, K. Carley, and P. Pattison, eds., *Dynamic Social Network Modeling and Analysis*, pp. 39–77, National Academies Press, Washington, DC (2003).
- [191] Freeman, L. C., *The Development of Social Network Analysis*, Empirical Press, Vancouver (2004).
- [192] Freeman, L. C., Borgatti, S. P., and White, D. R., Centrality in valued graphs: A measure of betweenness based on network flow, *Soc. Networks* **13**, 141–154 (1991).
- [193] Freeman, L. C., Freeman, S. C., and Michaelson, A. G., On human social intelligence, *J. Social Biol. Struct.* **11**, 415–425 (1988).
- [194] Freeman, L. C., Freeman, S. C., and Michaelson, A. G., How humans see social groups: A test of the Sailer–Gaulin models, *J. Quant. Anthropol.* **1**, 229–238 (1989).
- [195] Fronczak, A., Hołyst, J. A., Jedynek, M., and Sienkiewicz, J., Higher order clustering coefficients in Barabási–Albert networks, *Physica A* **316**, 688–694 (2002).
- [196] Fu, T. Z. J., Song, Q., and Chiu, D. M., The academic social networks, *Scientometrics* **101**, 203–239 (2014).
- [197] Galaskiewicz, J., *Social Organization of an Urban Grants Economy*, Academic Press, New York (1985).
- [198] Gallotti, R. and Barthelemy, M., The multilayer temporal network of public transport in Great Britain, *Sci. Data* **2**, 140056 (2015).
- [199] Gao, Y., Zheng, Z., and Qin, F., Analysis of Linux kernel as a complex network, *Chaos, Solitons & Fractals* **69**, 246–252 (2014).
- [200] Garfield, E., Citation indexes for science, *Science* **122**, 108–111 (1955).
- [201] Gastner, M. T., Spatial distributions: Density-equalizing map projections, facility location, and two-dimensional networks, Ph.D. thesis, University of Michigan (2005).
- [202] Gastner, M. T. and Newman, M. E. J., Optimal design of spatial distribution networks, *Phys. Rev. E* **74**, 016117 (2006).
- [203] Gastner, M. T. and Newman, M. E. J., The spatial structure of networks, *Eur. Phys. J. B* **49**, 247–252 (2006).
- [204] Girvan, M. and Newman, M. E. J., Community structure in social and biological networks, *Proc. Natl. Acad. Sci. USA* **99**, 7821–7826 (2002).
- [205] Gleich, D. F., Pagerank beyond the web, *SIAM Rev.* **57**, 321–363 (2015).
- [206] Gleiser, P. and Danon, L., Community structure in jazz, *Adv. Complex Syst.* **6**, 565–573 (2003).
- [207] Gleiss, P. M., Stadler, P. F., Wagner, A., and Fell, D. A., Relevant cycles in chemical reaction networks, *Adv. Complex Syst.* **4**, 207–226 (2001).
- [208] Golbeck, J., Grimes, J. M., and Rogers, A., Twitter use by the US Congress, *J. Am. Soc. Inform. Sci. Technol.* **61**, 1612–1621 (2010).
- [209] Goldstein, M. L., Morris, S. A., and Yen, G. G., Problems with fitting to the power-law distribution, *Eur. Phys. J. B* **41**, 255–258 (2004).
- [210] Goltsev, A. V., Dorogovtsev, S. N., and Mendes, J. F. F., K-core (bootstrap) percolation on complex networks: Critical phenomena and nonlocal effects, *Phys. Rev. E* **73**, 056101 (2006).
- [211] Goltsev, A. V., Dorogovtsev, S. N., and Mendes, J. F. F., Percolation on correlated networks, *Phys. Rev. E* **78**, 051105 (2008).
- [212] Goncalves, B., Perra, N., and Vespignani, A., Modeling users’ activity on Twitter networks: Validation of Dunbar’s number, *PLOS One* **6**, e22656 (2011).
- [213] Good, B. H., de Montjoye, Y.-A., and Clauset, A., Performance of modularity maximization in practical contexts, *Phys. Rev. E* **81**, 046106 (2010).
- [214] Grant, T. R., Dominance and association among members of a captive and a free-ranging group of grey kangaroos (*Macropus giganteus*), *Animal Behav.* **21**, 449–456 (1973).
- [215] Grassberger, P., On the critical behavior of the general epidemic process and dynamical percolation, *Math. Biosci.* **63**, 157–172 (1983).
- [216] Gress, B., Properties of the USPTO patent citation network: 1963–2002, *World Patent Inform.* **32**, 3–21 (2010).
- [217] Grindrod, P. and Higham, D. J., Evolving graphs: Dynamical models, inverse problems and propagation, *Proc. R. Soc. London A* **466**, 753–770 (2010).

REFERENCES

- [218] Grossman, J. W., The evolution of the mathematical research collaboration graph, *Congr. Numer.* **158**, 202–212 (2002).
- [219] Grossman, J. W. and Ion, P. D. F., On a portion of the well-known collaboration graph, *Congr. Numer.* **108**, 129–131 (1995).
- [220] Grujić, J., Movies recommendation networks as bipartite graphs, in M. Bubak, G. D. Albada, J. Dongarra, and P. M. A. Sloot, eds., *Proceedings of the 8th International Conference on Computational Science*, no. 5102 in Lecture Notes in Computer Science, pp. 576–583, Springer, Berlin (2008).
- [221] Guare, J., *Six Degrees of Separation: A Play*, Vintage, New York (1990).
- [222] Guilbeault, D., Becker, J., and Centola, D., Complex contagions: A decade in review, in S. Lehmann and Y. Ahn, eds., *Spreading Dynamics in Social Systems*, Springer Nature, Berlin (2018).
- [223] Guimerà, R., Danon, L., Díaz-Guilera, A., Giralt, F., and Arenas, A., Self-similar community structure in a network of human interactions, *Phys. Rev. E* **68**, 065103 (2003).
- [224] Guimerà, R., Mossa, S., Turtschi, A., and Amaral, L. A. N., The worldwide air transportation network: Anomalous centrality, community structure, and cities’ global roles, *PNAS* **102**, 7794–7799 (2005).
- [225] Guimerà, R. and Sales-Pardo, M., Missing and spurious interactions and the reconstruction of complex networks, *Proc. Natl. Acad. Sci. USA* **106**, 22073–22078 (2009).
- [226] Guimerà, R. and Sales-Pardo, M., A network inference method for large-scale unsupervised identification of novel drug-drug interactions, *PLOS Comput. Biol.* **9**, e1003374 (2013).
- [227] Guimerà, R., Sales-Pardo, M., and Amaral, L. A. N., Modularity from fluctuations in random graphs and complex networks, *Phys. Rev. E* **70**, 025101 (2004).
- [228] Gutenberg, B. and Richter, R. F., Frequency of earthquakes in California, *Bulletin of the Seismological Society of America* **34**, 185–188 (1944).
- [229] Harary, F., On the notion of balance of a signed graph, *Michigan Math. J.* **2**, 143–146 (1953).
- [230] Harary, F., *Graph Theory*, Perseus, Cambridge, MA (1995).
- [231] Helmstaedter, M., Cellular-resolution connectomics: Challenges of dense neural circuit reconstruction, *Nat. Methods* **10**, 501–507 (2013).
- [232] Hethcote, H. W., The mathematics of infectious diseases, *SIAM Rev.* **42**, 599–653 (2000).
- [233] Hidalgo, C. A. and Rodriguez-Sickert, C., The dynamics of a mobile phone network, *Physica A* **387**, 3017–3024 (2008).
- [234] Hines, P., Cotilla-Sanchez, E., and Blumsack, S., Do topological models provide good information about electricity infrastructure vulnerability?, *Chaos* **20**, 033122 (2010).
- [235] Hoff, P. D., Raftery, A. E., and Handcock, M. S., Latent space approaches to social network analysis, *J. Amer. Stat. Assoc.* **97**, 1090–1098 (2002).
- [236] Hofmann, T., Unsupervised learning by probabilistic latent semantic analysis, *Mach. Learn.* **42**, 177–196 (2001).
- [237] Holme, P., Core-periphery organization of complex networks, *Phys. Rev. E* **72**, 046111 (2005).
- [238] Holme, P., Edling, C. R., and Liljeros, F., Structure and time-evolution of an Internet dating community, *Soc. Networks* **26**, 155–174 (2004).
- [239] Holme, P. and Saramäki, J., Temporal networks, *Phys. Rep.* **519**, 97–125 (2012).
- [240] Hopkins, A. L., Network pharmacology: The next paradigm in drug discovery, *Nat. Chem. Biol.* **4**, 682–690 (2008).
- [241] Hou, H., Kretschmer, H., and Liu, Z., The structure of scientific collaboration networks in Scientometrics, *Scientometrics* **75**, 189–202 (2008).
- [242] Hric, D., Kaski, K., and Kivelä, M., Stochastic block model reveals the map of citation patterns and their evolution in time, preprint arxiv:1705.00018 (2017).
- [243] Hua, Y. and Zhu, D., Empirical analysis of the worldwide maritime transportation network, *Physica A* **388**, 2061–2071 (2009).
- [244] Huberman, B. A., *The Laws of the Web*, MIT Press, Cambridge, MA (2001).
- [245] Huxham, M., Beaney, S., and Raffaelli, D., Do parasites reduce the chances of triangulation in a real food web?, *Oikos* **76**, 284–300 (1996).
- [246] Jacobs, A. Z., Way, S. F., Ugander, J., and Clauset, A., Assembling thefacebook: Using heterogeneity to understand online social network assembly, in

- Proceedings of the 7th Annual ACM Web Science Conference*, p. 18, Association of Computing Machinery, New York (2015).
- [247] Jaffe, A. and Trajtenberg, M., *Patents, Citations and Innovations: A Window on the Knowledge Economy*, MIT Press, Cambridge, MA (2002).
- [248] Jeh, G. and Widom, J., SimRank: A measure of structural-context similarity, in *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 538–543, Association of Computing Machinery, New York (2002).
- [249] Jenks, S. M. and Ginsburg, B. E., Socio-sexual dynamics in a captive wolf pack, in H. Frank, ed., *Man and Wolf*, pp. 375–399, Junk Publishers, Dordrecht (1987).
- [250] Jeong, H., Mason, S., Barabási, A.-L., and Oltvai, Z. N., Lethality and centrality in protein networks, *Nature* **411**, 41–42 (2001).
- [251] Jeong, H., Néda, Z., and Barabási, A.-L., Measuring preferential attachment in evolving networks, *Europhys. Lett.* **61**, 567–572 (2003).
- [252] Jeong, H., Tombor, B., Albert, R., Oltvai, Z. N., and Barabási, A.-L., The large-scale organization of metabolic networks, *Nature* **407**, 651–654 (2000).
- [253] Jones, J. H. and Handcock, M. S., Sexual contacts and epidemic thresholds, *Nature* **423**, 605–606 (2003).
- [254] Kansky, K. J., *Structure of Transportation Networks: Relationships Between Network Geometry and Regional Characteristics*, University of Chicago, Chicago (1963).
- [255] Karrer, B. and Newman, M. E. J., Random graph models for directed acyclic networks, *Phys. Rev. E* **80**, 046110 (2009).
- [256] Karrer, B. and Newman, M. E. J., Random graphs containing arbitrary distributions of subgraphs, *Phys. Rev. E* **82**, 066118 (2010).
- [257] Karrer, B. and Newman, M. E. J., Stochastic block-models and community structure in networks, *Phys. Rev. E* **83**, 016107 (2011).
- [258] Katz, L., A new status index derived from sociometric analysis, *Psychometrika* **18**, 39–43 (1953).
- [259] Killworth, P. D. and Bernard, H. R., Informant accuracy in social network data, *Hum. Organ.* **35**, 269–286 (1976).
- [260] Killworth, P. D. and Bernard, H. R., The reverse small world experiment, *Soc. Networks* **1**, 159–192 (1978).
- [261] Killworth, P. D., Johnsen, E. C., Bernard, H. R., Shelley, G. A., and McCarty, C., Estimating the size of personal networks, *Soc. Networks* **12**, 289–312 (1990).
- [262] Kim, J., Krapivsky, P. L., Kahng, B., and Redner, S., Infinite-order percolation and giant fluctuations in a protein interaction network, *Phys. Rev. E* **66**, 055101 (2002).
- [263] Kinney, R., Crucitti, P., Albert, R., and Latora, V., Modeling cascading failures in the North American power grid, *Eur. Phys. J. B* **46**, 101–107 (2004).
- [264] Kivelä, M., Arenas, A., Barthelemy, M., Gleeson, J. P., Moreno, Y., and Porter, M. A., Multilayer networks, *J. Complex Netw.* **2**, 203–271 (2014).
- [265] Kleinberg, J. M., Authoritative sources in a hyperlinked environment, *J. ACM* **46**, 604–632 (1999).
- [266] Kleinberg, J. M., Navigation in a small world, *Nature* **406**, 845 (2000).
- [267] Kleinberg, J. M., The small-world phenomenon: An algorithmic perspective, in *Proceedings of the 32nd Annual ACM Symposium on Theory of Computing*, pp. 163–170, Association of Computing Machinery, New York (2000).
- [268] Kleinberg, J. M., Small world phenomena and the dynamics of information, in T. G. Dietterich, S. Becker, and Z. Ghahramani, eds., *Proceedings of the 2001 Neural Information Processing Systems Conference*, MIT Press, Cambridge, MA (2002).
- [269] Kleinberg, J. M., Kumar, S. R., Raghavan, P., Rajagopalan, S., and Tomkins, A., The Web as a graph: Measurements, models and methods, in T. Asano, H. Imai, D. T. Lee, S.-I. Nakano, and T. Tokuyama, eds., *Proceedings of the 5th Annual International Conference on Combinatorics and Computing*, no. 1627 in Lecture Notes in Computer Science, pp. 1–18, Springer, Berlin (1999).
- [270] Klovdahl, A. S., Urban social networks: Some methodological problems and possibilities, in M. Kochen, ed., *The Small World*, Ablex Publishing, Norwood, NJ (1989).
- [271] Klovdahl, A. S., Potterat, J. J., Woodhouse, D. E., Muth, J. B., Muth, S. Q., and Darrow, W. W., Social

REFERENCES

- networks and infectious disease: The Colorado Springs study, *Soc. Sci. Med.* **38**, 79–88 (1994).
- [272] Knuth, D. E., *The Stanford GraphBase: A Platform for Combinatorial Computing*, Addison-Wesley, Reading, MA (1993).
- [273] Kohli, R. and Sah, R., Some empirical regularities in market shares, *Management Science* **52**, 1792–1798 (2006).
- [274] Koren, Y., Drawing graphs by eigenvectors: Theory and practice, *Comput. Math. Appl.* **49**, 1867–1888 (2005).
- [275] Korte, C. and Milgram, S., Acquaintance links between racial groups: Application of the small world method, *J. Pers. Soc. Psychol.* **15**, 101–108 (1970).
- [276] Kossinets, G., Effects of missing data in social networks, *Soc. Networks* **28**, 247–268 (2006).
- [277] Kossinets, G. and Watts, D. J., Empirical analysis of an evolving social network, *Science* **311**, 88–90 (2006).
- [278] Kramer, A. D. I., Guillory, J. E., and Hancock, J. T., Experimental evidence of massive-scale emotional contagion through social networks, *Proc. Natl. Acad. Sci. USA* **111**, 8788–8790 (2014).
- [279] Krapivsky, P. L. and Redner, S., Organization of growing random networks, *Phys. Rev. E* **63**, 066123 (2001).
- [280] Krapivsky, P. L., Redner, S., and Leyvraz, F., Connectivity of growing random networks, *Phys. Rev. Lett.* **85**, 4629–4632 (2000).
- [281] Krapivsky, P. L., Rodgers, G. J., and Redner, S., Degree distributions of growing networks, *Phys. Rev. Lett.* **86**, 5401–5404 (2001).
- [282] Krebs, V. E., Mapping networks of terrorist cells, *Connections* **24**, 43–52 (2002).
- [283] Krogan, N. J., Cagney, G., Yu, H., Zhong, G., Guo, X., Ignatchenko, A., Li, J., Pu, S., Datta, N., Tikuisis, A. P., Punna, T., Peregrín-Alvarez, J. M., Shales, M., Zhang, X., Davey, M., Robinson, M. D., Paccanaro, A., Bray, J. E., Sheung, A., Beattie, B., Richards, D. P., Canadien, V., Lalev, A., Mena, F., Wong, P., Starostine, A., Canete, M. M., Vlasblom, J., Wu, S., Orsi, C., Collins, S. R., Chandran, S., Haw, R., Rilstone, J. J., Gandi, K., Thompson, N. J., Musso, G., Onge, P. S., Ghanny, S., Lam, M. H. Y., Butland, G., Altaf-Ul, A. M., Kanaya, S., Shilatifard, A., O’Shea, E., Weissman, J. S., Ingles, C. J., Hughes, T. R., Parkinson, J., Gerstein, M., Wodak, S. J., Emili, A., and Greenblatt, J. F., Global landscape of protein complexes in the yeast *Saccharomyces cerevisiae*, *Nature* **440**, 637–643 (2006).
- [284] Lakhina, A., Byers, J., Crovella, M., and Xie, P., Sampling biases in IP topology measurements, in *Proceedings of the 22nd Annual Joint Conference of the IEEE Computer and Communications Societies*, Institute of Electrical and Electronics Engineers, New York (2003).
- [285] Lambiotte, R., Geographical dispersal of mobile communication networks, *Physica A* **387**, 5317–5325 (2008).
- [286] Lancichinetti, A. and Fortunato, S., Community detection algorithms: A comparative analysis, *Phys. Rev. E* **80**, 056117 (2009).
- [287] Lancichinetti, A., Fortunato, S., and Radicchi, F., Benchmark graphs for testing community detection algorithms, *Phys. Rev. E* **78**, 046110 (2008).
- [288] Landauer, T. K., Foltz, P. W., and Laham, D., An introduction to latent semantic analysis, *Discourse Process.* **25**, 259–284 (1998).
- [289] Langville, A. N. and Meyer, C. D., *Who’s #1? The Science of Rating and Ranking*, Princeton University Press, Princeton, NJ (2013).
- [290] Latora, V. and Marchiori, M., Is the Boston subway a small-world network?, *Physica A* **314**, 109–113 (2002).
- [291] Lee, D. D. and Seung, H. S., Learning the parts of objects by nonnegative matrix factorization, *Nature* **401**, 788–791 (1999).
- [292] Lee, D. D. and Seung, H. S., Algorithms for non-negative matrix factorization, in *Proceedings of the 2000 Neural Information Processing Systems Conference*, pp. 556–562, MIT Press, Cambridge, MA (2001).
- [293] Lee, K., Jung, W.-S., Park, J. S., and Choi, M. Y., Statistical analysis of the Metropolitan Seoul Subway system: Network structure and passenger flows, *Physica A* **387**, 6231–6234 (2008).
- [294] Lehmann, S., Lautrup, B., and Jackson, A. D., Citation networks in high energy physics, *Phys. Rev. E* **68**, 026113 (2003).

- [295] Leicht, E. A., Clarkson, G., Shedden, K., and Newman, M. E. J., Large-scale structure of time evolving citation networks, *Eur. Phys. J. B* **59**, 75–83 (2007).
- [296] Leicht, E. A., Holme, P., and Newman, M. E. J., Vertex similarity in networks, *Phys. Rev. E* **73**, 026120 (2006).
- [297] Lewis, K., The limits of racial prejudice, *Proc. Natl. Acad. Sci. USA* **110**, 18814–18819 (2013).
- [298] Lewis, K., Kaufman, J., Gonzalez, M., Wimmer, A., and Christakis, N., Tastes, ties, and time: A new social network dataset using Facebook.com, *Soc. Networks* **30**, 330–342 (2008).
- [299] Li, S., Chen, Y., Du, H., and Feldman, M. W., A genetic algorithm with local search strategy for improved detection of community structure, *Complexity* **15**(4), 53–60 (2010).
- [300] Liben-Nowell, D., Geographic routing in social networks, *Proc. Natl. Acad. Sci. USA* **102**, 11623–11628 (2005).
- [301] Liben-Nowell, D. and Kleinberg, J., The link-prediction problem for social networks, *J. Assoc. Inf. Sci. Technol.* **58**, 1019–1031 (2007).
- [302] Lichtman, J. W., Livet, J., and Sanes, J. R., A technicolour approach to the connectome, *Nat. Rev. Neurosci.* **9**, 417–422 (2008).
- [303] Liggett, T. M., *Interacting Particle Systems*, Springer, New York (1985).
- [304] Liljeros, F., Edling, C. R., and Amaral, L. A. N., Sexual networks: Implication for the transmission of sexually transmitted infection, *Microbes Infect.* **5**, 189–196 (2003).
- [305] Liljeros, F., Edling, C. R., Amaral, L. A. N., Stanley, H. E., and Åberg, Y., The web of human sexual contacts, *Nature* **411**, 907–908 (2001).
- [306] Lloyd, A. L. and May, R. M., How viruses spread among computers and people, *Science* **292**, 1316–1317 (2001).
- [307] Lloyd, J. R., Orbanz, P., Ghahramani, Z. and Roy, D. M., Random function priors for exchangeable arrays with applications to graphs and relational data, in *Proceedings of the 2012 Neural Information Processing Systems Conference*, pp. 1–9, MIT Press, Cambridge, MA (2012).
- [308] Lorenz, M. O., Methods of measuring the concentration of wealth, *Publ. Am. Stat. Assoc.* **9**, 209–219 (1905).
- [309] Lotka, A. J., The frequency distribution of scientific production, *J. Wash. Acad. Sci.* **16**, 317–323 (1926).
- [310] Lott, D. F., Dominance relations and breeding rate in mature male American bison, *Z. Tierpsychol.* **49**, 418–432 (1979).
- [311] Lowry, O. H., Rosebrough, N. J., Farr, A. L., and Randall, R. J., Protein measurement with the Folin phenol reagent, *J. Biol. Chem.* **193**, 265–275 (1951).
- [312] Lu, E. T. and Hamilton, R. J., Avalanches of the distribution of solar flares, *Astrophys. J.* **380**, 89–92 (1991).
- [313] Lueg, C. and Fisher, D., eds., *From Usenet to CoWebs: Interacting with Social Information Spaces*, Springer, New York (2003).
- [314] Lupu, Y. and Voeten, E., Precedent in international courts: A network analysis of case citations by the European Court of Human Rights, *Br. J. Polit. Sci.* **42**, 413–439 (2012).
- [315] Lusseau, D., The emergent properties of a dolphin social network, *Proc. R. Soc. London B (suppl.)* **270**, S186–S188 (2003).
- [316] Lusseau, D., Schneider, K., Boisseau, O. J., Haase, P., Slooten, E., and Dawson, S. M., The bottlenose dolphin community of Doubtful Sound features a large proportion of long-lasting associations. Can geographic isolation explain this unique trait?, *Behav. Ecol. Sociobiol.* **54**, 396–405 (2003).
- [317] MacKinnon, I. and Warren, R., Age and geographic inferences of the LiveJournal social network, in E. Airolidi, D. M. Blei, S. E. Fienberg, A. Goldenberg, E. P. Xing, and A. X. Zheng, eds., *Statistical Network Analysis: Models, Issues, and New Directions*, vol. 4503 of *Lecture Notes in Computer Science*, pp. 176–178, Springer-Verlag, Berlin (2007).
- [318] Mariolis, P., Interlocking directorates and control of corporations: The theory of bank control, *Soc. Sci. Q.* **56**, 425–439 (1975).
- [319] Maritan, A., Rinaldo, A., Rigon, R., Giacometti, A., and Rodríguez-Iturbe, I., Scaling laws for river networks, *Phys. Rev. E* **53**, 1510–1515 (1996).
- [320] Marsden, P. V., Network data and measurement, *Annu. Rev. Sociol.* **16**, 435–463 (1990).
- [321] Martinez, N. D., Artifacts or attributes? Effects of resolution on the Little Rock Lake food web, *Ecol. Monographs* **61**, 367–392 (1991).

REFERENCES

- [322] Martinez, N. D., Constant connectance in community food webs, *Am. Natur.* **139**, 1208–1218 (1992).
- [323] Maslov, S., Sneppen, K., and Zaliznyak, A., Detection of topological patterns in complex networks: Correlation profile of the internet, *Physica A* **333**, 529–540 (2004).
- [324] Masucci, A. P., Smith, D., Crooks, A., and Batty, M., Random planar graphs and the London street network, *Eur. Phys. J. B* **71**, 259–271 (2009).
- [325] May, R. M. and Anderson, R. M., The transmission dynamics of human immunodeficiency virus (HIV), *Philos. Trans. R. Soc. London B* **321**, 565–607 (1988).
- [326] McCarty, C., Killworth, P. D., Bernard, H. R., Johnsen, E. C., and Shelley, G. A., Comparing two methods for estimating network size, *Hum. Organ.* **60**, 28–39 (2001).
- [327] McDaid, A. F., Greene, D., and Hurley, N., Normalized mutual information to evaluate overlapping community finding algorithms, preprint arxiv:1110.2515 (2011).
- [328] McMahan, C. A. and Morris, M. D., Application of maximum likelihood paired comparison ranking to estimation of a linear dominance hierarchy in animal societies, *Animal Behav.* **32**, 374–378 (1984).
- [329] Medus, A., Acuña, G., and Dorso, C. O., Detection of community structures in networks via global optimization, *Physica A* **358**, 593–604 (2005).
- [330] Menger, K., Zur allgemeinen Kurventheorie, *Fundamenta Mathematicae* **10**, 96–115 (1927).
- [331] Meyer, C. D., *Matrix Analysis and Applied Linear Algebra*, Society for Industrial and Applied Mathematics, Philadelphia (2000).
- [332] Mézard, M. and Montanari, A., *Information, Physics, and Computation*, Oxford University Press, Oxford (2009).
- [333] Milgram, S., The small world problem, *Psychol. Today* **2**, 60–67 (1967).
- [334] Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D., and Alon, U., Network motifs: Simple building blocks of complex networks, *Science* **298**, 824–827 (2002).
- [335] Mitzenmacher, M., A brief history of generative models for power law and lognormal distributions, *Internet Math.* **1**, 226–251 (2004).
- [336] Mollison, D., Spatial contact models for ecological and epidemic spread, *J. R. Stat. Soc. B* **39**, 283–326 (1977).
- [337] Molloy, M. and Reed, B., A critical point for random graphs with a given degree sequence, *Random Struct. Alg.* **6**, 161–179 (1995).
- [338] Moody, J., Race, school integration, and friendship segregation in America, *Am. J. Sociol.* **107**, 679–716 (2001).
- [339] Moore, C., Ghoshal, G., and Newman, M. E. J., Exact solutions for models of evolving networks with addition and deletion of nodes, *Phys. Rev. E* **74**, 036121 (2006).
- [340] Moore, C. and Mertens, S., *The Nature of Computation*, Oxford University Press, Oxford (2011).
- [341] Moreno, J. L., *Who Shall Survive?*, Beacon House, Beacon, NY (1934).
- [342] Moreno, Y., Pastor-Satorras, R., and Vespignani, A., Epidemic outbreaks in complex heterogeneous networks, *Eur. Phys. J. B* **26**, 521–529 (2002).
- [343] Myers, C. R., Software systems as complex networks: Structure, function, and evolvability of software collaboration graphs, *Phys. Rev. E* **68**, 046116 (2003).
- [344] Neukum, G. and Ivanov, B. A., Crater size distributions and impact probabilities on Earth from lunar, terrestrial-planet, and asteroid cratering data, in T. Gehrels, ed., *Hazards Due to Comets and Asteroids*, pp. 359–416, University of Arizona Press, Tucson, AZ (1994).
- [345] Newman, E. I., A method of estimating the total length of root in a sample, *J. Appl. Ecol.* **3**, 139–145 (1966).
- [346] Newman, M. E. J., Clustering and preferential attachment in growing networks, *Phys. Rev. E* **64**, 025102 (2001).
- [347] Newman, M. E. J., Scientific collaboration networks: I. Network construction and fundamental results, *Phys. Rev. E* **64**, 016131 (2001).
- [348] Newman, M. E. J., Scientific collaboration networks: II. Shortest paths, weighted networks, and centrality, *Phys. Rev. E* **64**, 016132 (2001).
- [349] Newman, M. E. J., The structure of scientific collaboration networks, *Proc. Natl. Acad. Sci. USA* **98**, 404–409 (2001).

- [350] Newman, M. E. J., Assortative mixing in networks, *Phys. Rev. Lett.* **89**, 208701 (2002).
- [351] Newman, M. E. J., Ego-centered networks and the ripple effect, *Soc. Networks* **25**, 83–95 (2003).
- [352] Newman, M. E. J., Mixing patterns in networks, *Phys. Rev. E* **67**, 026126 (2003).
- [353] Newman, M. E. J., Properties of highly clustered networks, *Phys. Rev. E* **68**, 026121 (2003).
- [354] Newman, M. E. J., The structure and function of complex networks, *SIAM Rev.* **45**, 167–256 (2003).
- [355] Newman, M. E. J., Analysis of weighted networks, *Phys. Rev. E* **70**, 056131 (2004).
- [356] Newman, M. E. J., A measure of betweenness centrality based on random walks, *Soc. Networks* **27**, 39–54 (2005).
- [357] Newman, M. E. J., Power laws, Pareto distributions and Zipf's law, *Contemp. Phys.* **46**, 323–351 (2005).
- [358] Newman, M. E. J., Threshold effects for two pathogens spreading on a network, *Phys. Rev. Lett.* **95**, 108701 (2005).
- [359] Newman, M. E. J., Modularity and community structure in networks, *Proc. Natl. Acad. Sci. USA* **103**, 8577–8582 (2006).
- [360] Newman, M. E. J., Random graphs with clustering, *Phys. Rev. Lett.* **103**, 058701 (2009).
- [361] Newman, M. E. J., Network reconstruction and error estimation with noisy network data, preprint arxiv:1803.02427 (2018).
- [362] Newman, M. E. J., Network structure from rich but noisy data, *Nat. Phys.* **14**, 542–545 (2018).
- [363] Newman, M. E. J. and Ferrario, C. R., Interacting epidemics and coinfection on contact networks, *PLOS One* **8**, e71321 (2013).
- [364] Newman, M. E. J., Forrest, S., and Balthrop, J., Email networks and the spread of computer viruses, *Phys. Rev. E* **66**, 035101 (2002).
- [365] Newman, M. E. J. and Girvan, M., Mixing patterns and community structure in networks, in R. Pastor-Satorras, J. Rubi, and A. Díaz-Guilera, eds., *Statistical Mechanics of Complex Networks*, no. 625 in Lecture Notes in Physics, pp. 66–87, Springer, Berlin (2003).
- [366] Newman, M. E. J. and Girvan, M., Finding and evaluating community structure in networks, *Phys. Rev. E* **69**, 026113 (2004).
- [367] Newman, M. E. J. and Park, J., Why social networks are different from other types of networks, *Phys. Rev. E* **68**, 036122 (2003).
- [368] Newman, M. E. J. and Peixoto, T. P., Generalized communities in networks, *Phys. Rev. Lett.* **115**, 088701 (2015).
- [369] Newman, M. E. J., Strogatz, S. H., and Watts, D. J., Random graphs with arbitrary degree distributions and their applications, *Phys. Rev. E* **64**, 026118 (2001).
- [370] Newman, M. E. J. and Watts, D. J., Scaling and percolation in the small-world network model, *Phys. Rev. E* **60**, 7332–7342 (1999).
- [371] Newman, M. E. J. and Ziff, R. M., Fast Monte Carlo algorithm for site or bond percolation, *Phys. Rev. E* **64**, 016706 (2001).
- [372] Ng, A. Y., Zheng, A. X., and Jordan, M. I., Stable algorithms for link analysis, in D. H. Kraft, W. B. Croft, D. J. Harper, and J. Zobel, eds., *Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 258–266, Association of Computing Machinery, New York (2001).
- [373] Ogielski, A. T., Integer optimization and zero-temperature fixed point in Ising random-field systems, *Phys. Rev. Lett.* **57**, 1251–1254 (1986).
- [374] Onnela, J.-P., Arbesman, S., González M. C., Barabási, A.-L., and Christakis, N. A., Geographic constraints on social network groups, *PLOS One* **6**, e16939 (2011).
- [375] Onnela, J.-P., Saramäki, J., Hyvönen, J., Szabó, G., Lazer, D., Kaski, K., Kertész, J., and Barabási, A.-L., Structure and tie strengths in mobile communication networks, *Proc. Natl. Acad. Sci. USA* **104**, 7332–7336 (2007).
- [376] Orman, G. K., Labatut, V., and Cherifi, H., Qualitative comparison of community detection algorithms, *Commun. Computer Inform. Sci.* **167**, 265–279 (2011).
- [377] Padgett, J. F. and Ansell, C. K., Robust action and the rise of the Medici, 1400–1434, *Am. J. Sociol.* **98**, 1259–1319 (1993).
- [378] Pagani, G. A. and Aiello, M., The power grid as a complex network: A survey, *Physica A* **392**, 2688–2700 (2013).

REFERENCES

- [379] Palla, G., Derényi, I., Farkas, I., and Vicsek, T., Uncovering the overlapping community structure of complex networks in nature and society, *Nature* **435**, 814–818 (2005).
- [380] Park, J., Diagrammatic perturbation methods in networks and sports ranking combinatorics, *J. Stat. Mech.* **2010**, P04006 (2010).
- [381] Pastor-Satorras, R., Vázquez, A., and Vespignani, A., Dynamical and correlation properties of the Internet, *Phys. Rev. Lett.* **87**, 258701 (2001).
- [382] Pastor-Satorras, R. and Vespignani, A., Epidemic dynamics and endemic states in complex networks, *Phys. Rev. E* **63**, 066117 (2001).
- [383] Pastor-Satorras, R. and Vespignani, A., Epidemic spreading in scale-free networks, *Phys. Rev. Lett.* **86**, 3200–3203 (2001).
- [384] Pastor-Satorras, R. and Vespignani, A., *Evolution and Structure of the Internet*, Cambridge University Press, Cambridge (2004).
- [385] Peixoto, T. P., Hierarchical block structures and high-resolution model selection in large networks, *Phys. Rev. X* **4**, 011047 (2014).
- [386] Pelletier, J. D., Self-organization and scaling relationships of evolving river networks, *J. Geophys. Res.* **104**, 7359–7375 (1999).
- [387] Pitts, F. R., A graph theoretic approach to historical geography, *The Professional Geographer* **17**, 15–20 (1965).
- [388] Plischke, M. and Bergersen, B., *Equilibrium Statistical Physics*, 3rd edn., World Scientific, Singapore (2006).
- [389] Pool, I. de S. and Kochen, M., Contacts and influence, *Soc. Networks* **1**, 1–48 (1978).
- [390] Porter, M. A. and Gleeson, J., *Dynamical Systems on Networks: A Tutorial*, Springer, Berlin (2016).
- [391] Pothen, A., Simon, H., and Liou, K.-P., Partitioning sparse matrices with eigenvectors of graphs, *SIAM J. Matrix Anal. Appl.* **11**, 430–452 (1990).
- [392] Potterat, J. J., Phillips-Plummer, L., Muth, S. Q., Rothenberg, R. B., Woodhouse, D. E., Maldonado-Long, T. S., Zimmerman, H. P., and Muth, J. B., Risk network structure in the early epidemic phase of HIV transmission in Colorado Springs, *Sex. Transm. Infect.* **78**, i159–i163 (2002).
- [393] Price, D. J. de S., Networks of scientific papers, *Science* **149**, 510–515 (1965).
- [394] Price, D. J. de S., A general theory of bibliometric and other cumulative advantage processes, *J. Amer. Soc. Inform. Sci.* **27**, 292–306 (1976).
- [395] Radicchi, F., Castellano, C., Cecconi, F., Loreto, V., and Parisi, D., Defining and identifying communities in networks, *Proc. Natl. Acad. Sci. USA* **101**, 2658–2663 (2004).
- [396] Radicchi, F., Fortunato, S., and Castellano, C., Universality of citation distributions: Towards an objective measure of scientific impact, *Proc. Natl. Acad. Sci. USA* **105**, 17268–17272 (2008).
- [397] Radicchi, F., Fortunato, S., Markines, B., and Vespignani, A., Diffusion of scientific credits and the ranking of scientists, *Phys. Rev. E* **80**, 056103 (2009).
- [398] Radicchi, F., Fortunato, S., and Vespignani, A., Citation networks, in A. Scharnhorst, K. Börner, and P. van den Besselaar, eds., *Models of Science Dynamics: Encounters Between Complexity Theory and Information Sciences*, pp. 233–257, Springer, New York (2012).
- [399] Rand, W. M., Objective criteria for the evaluations of clustering methods, *J. Am. Stat. Assoc.* **66**, 846–850 (1971).
- [400] Rapoport, A. and Horvath, W. J., A study of a large sociogram, *Behav. Sci.* **6**, 279–291 (1961).
- [401] Ratti, C., Sobolevsky, S., Calabrese, F., Andris, C., Reades, J., Martino, M., Claxton, R., and Strogatz, S. H., Redrawing the map of Great Britain from a network of human interactions, *PLOS One* **5**, e14248 (2010).
- [402] Ravasz, E. and Barabási, A.-L., Hierarchical organization in complex networks, *Phys. Rev. E* **67**, 026112 (2003).
- [403] Rea, L. M. and Parker, R. A., *Designing and Conducting Survey Research: A Comprehensive Guide*, 3rd edn., Jossey-Bass, San Francisco, CA (2005).
- [404] Redner, S., How popular is your paper? An empirical study of the citation distribution, *Eur. Phys. J. B* **4**, 131–134 (1998).
- [405] Redner, S., Citation statistics from 110 years of Physical Review, *Physics Today* **58**, 49–54 (2005).
- [406] Ricci, F., Rokach, L., and Shapira, B., eds., *Recommender Systems Handbook*, 2nd edn., Springer, Berlin (2015).

- [407] Rinaldo, A., Banavar, J. R., and Maritan, A., Trees, networks, and hydrology, *Water Resour. Res.* **42**, W06D07 (2006).
- [408] Riolo, M. A., Cantwell, G. T., Reinert, G., and Newman, M. E. J., Efficient method for estimating the number of communities in a network, *Phys. Rev. E* **96**, 032310 (2017).
- [409] Ripeanu, M., Foster, I., and Iamnitchi, A., Mapping the Gnutella network: Properties of large-scale peer-to-peer systems and implications for system design, *IEEE Internet Comput.* **6**, 50–57 (2002).
- [410] Roberts, D. C. and Turcotte, D. L., Fractality and self-organized criticality of wars, *Fractals* **6**, 351–357 (1998).
- [411] Rocha, L. E. C., Liljeros, F., and Holme, P., Information dynamics shape the sexual networks of Internet-mediated prostitution, *Proc. Natl. Acad. Sci. USA* **107**, 5706–5711 (2010).
- [412] Rodríguez-Iturbe, I. and Rinaldo, A., *Fractal River Basins: Chance and Self-Organization*, Cambridge University Press, Cambridge (1997).
- [413] Rohe, K., Network driven sampling: A critical threshold for design effects, preprint arxiv:1505.05461 (2016).
- [414] Rombach, M. P., Porter, M. A., Fowler, J. H., and Mucha, P. J., Core-periphery structure in networks, *SIAM J. Appl. Math.* **74**, 167–190 (2014).
- [415] Rosas-Casals, M., Valverde, S., and Solé, R. V., Topological vulnerability of the European power grid under errors and attacks, *Int. J. Bifurcation Chaos* **17**, 2465–2475 (2007).
- [416] Rosvall, M. and Bergstrom, C. T., Maps of random walks on complex networks reveal community structure, *Proc. Natl. Acad. Sci. USA* **105**, 1118–1123 (2008).
- [417] Rothenberg, R., Baldwin, J., Trotter, R., and Muth, S., The risk environment for HIV transmission: Results from the Atlanta and Flagstaff network studies, *J. Urban Health* **78**, 419–431 (2001).
- [418] Sade, D. S., Sociometrics of Macaca mulatta: I. Linkages and cliques in grooming matrices, *Folia Primatol.* **18**, 196–223 (1972).
- [419] Sailer, L. D. and Gaulin, S. J. C., Proximity, sociality and observation: The definition of social groups, *Am. Anthropol.* **86**, 91–98 (1984).
- [420] Salganik, M. J., Dodds, P. S., and Watts, D. J., Experimental study of inequality and unpredictability in an artificial cultural market, *Science* **311**, 854–856 (2006).
- [421] Salganik, M. J. and Heckathorn, D. D., Sampling and estimation in hidden populations using respondent-driven sampling, *Sociol. Methodol.* **34**, 193–239 (2004).
- [422] Salton, G., *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*, Addison-Wesley, Reading, MA (1989).
- [423] Schank, T. and Wagner, D., Approximating clustering coefficient and transitivity, *J. Graph Algorithms Appl.* **9**, 265–275 (2005).
- [424] Scott, J., *Social Network Analysis: A Handbook*, 2nd edn., Sage, London (2000).
- [425] Sen, P., Dasgupta, S., Chatterjee, A., Sreeram, P. A., Mukherjee, G., and Manna, S. S., Small-world properties of the Indian railway network, *Phys. Rev. E* **67**, 036106 (2003).
- [426] Sienkiewicz, J. and Holyst, J. A., Statistical analysis of 22 public transport networks in Poland, *Phys. Rev. E* **72**, 046127 (2005).
- [427] Simkin, M. V. and Roychowdhury, V. P., Read before you cite, *Complex Syst.* **14**, 269–274 (2003).
- [428] Simkin, M. V. and Roychowdhury, V. P., Stochastic modeling of citation slips, *Scientometrics* **62**, 367–384 (2005).
- [429] Simon, H. A., On a class of skew distribution functions, *Biometrika* **42**, 425–440 (1955).
- [430] Smalheiser, N. R. and Torvik, V. I., Author name disambiguation, *Annu. Rev. Inform. Sci. Technol.* **43**, 287–313 (2009).
- [431] Smith, M., Invisible crowds in cyberspace: Measuring and mapping the social structure of USENET, in M. Smith and P. Kollock, eds., *Communities in Cyberspace*, Routledge Press, London (1999).
- [432] Solé, R. V., Pastor-Satorras, R., Smith, E., and Kepler, T. B., A model of large-scale proteome evolution, *Adv. Complex Syst.* **5**, 43–54 (2002).
- [433] Solomonoff, R. and Rapoport, A., Connectivity of random nets, *Bull. Math. Biophys.* **13**, 107–117 (1951).
- [434] Sood, V. and Redner, S., Voter model on heterogeneous graphs, *Phys. Rev. Lett.* **94**, 178701 (2005).

REFERENCES

- [435] Spielman, D. A. and Teng, S.-H., Spectral sparsification of graphs, *SIAM J. Comput.* **40**, 981–1025 (2011).
- [436] Sporns, O., *Networks of the Brain*, MIT Press, Cambridge, MA (2010).
- [437] Stauffer, D. and Aharony, A., *Introduction to Percolation Theory*, 2nd edn., Taylor and Francis, London (1992).
- [438] Stoica, I., Morris, R., Karger, D., Kaashoek, M. F., and Balakrishnan, H., Chord: A scalable peer-to-peer lookup service for Internet applications, in *Proceedings of the 2001 ACM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications (SIGCOMM)*, pp. 149–160, Association of Computing Machinery, New York (2001).
- [439] Stopczynski, A., Sekara, V., Sapiezynski, P., Cuttone, A., Madsen, M. M., Larsen, J. E., and Lehmann, S., Measuring large-scale social networks with high resolution, *PLOS One* **9**, e95978 (2014).
- [440] Strang, G., *Introduction to Linear Algebra*, Wellesley Cambridge Press, Wellesley, MA (2009).
- [441] Strogatz, S. H., *Nonlinear Dynamics and Chaos*, 2nd edn., Westview Press, Boulder (2014).
- [442] Stutzbach, D. and Rejaie, R., Characterizing today’s Gnutella topology, Technical Report CIS-TR-04-02, Department of Computer Science, University of Oregon (2004).
- [443] Szabó, G., Alava, M., and Kertész, J., Structural transitions in scale-free networks, *Phys. Rev. E* **67**, 056102 (2003).
- [444] Tang, J., Fong, A. C. M., Wang, B., and Zhang, J., A unified probabilistic framework for name disambiguation in digital library, *IEEE Trans. Knowl. Data Eng.* **24**, 975–987 (2012).
- [445] Thompson, S. K. and Frank, O., Model-based estimation with link-tracing sampling designs, *Surv. Methodol.* **26**, 87–98 (2000).
- [446] Traud, A. L., Mucha, P. J., and Porter, M. A., Social structure of Facebook networks, *Physica A* **391**, 4165–4180 (2012).
- [447] Travers, J. and Milgram, S., An experimental study of the small world problem, *Sociometry* **32**, 425–443 (1969).
- [448] Tsourakakis, C. E., Fast counting of triangles in large real networks without counting: Algorithms and laws, in *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pp. 608–617, Institute of Electrical and Electronics Engineers, New York (2008).
- [449] Turner, T. C., Smith, M. A., Fisher, D., and Welser, H. T., Picturing Usenet: Mapping computer-mediated collective action, *J. Comput.-Mediat. Commun.* **10**(4), 7 (2005).
- [450] Tyler, J. R., Wilkinson, D. M., and Huberman, B. A., Email as spectroscopy: Automated discovery of community structure within organizations, in M. Huysman, E. Wenger, and V. Wulf, eds., *Proceedings of the First International Conference on Communities and Technologies*, Kluwer, Dordrecht (2003).
- [451] Udry, J. R., Bearman, P. S., and Harris, K. M., National Longitudinal Study of Adolescent Health (1997). This research uses data from Add Health, a program project directed by Kathleen Mullan Harris and designed by J. Richard Udry, Peter S. Bearman, and Kathleen Mullan Harris at the University of North Carolina at Chapel Hill, and funded by grant P01–HD31921 from the Eunice Kennedy Shriver National Institute of Child Health and Human Development, with cooperative funding from 23 other federal agencies and foundations. Special acknowledgment is due Ronald R. Rindfuss and Barbara Entwisle for assistance in the original design. Information on how to obtain the Add Health data files is available on the Add Health website (<http://www.cpc.unc.edu/addhealth>). No direct support was received from grant P01–HD31921 for this analysis.
- [452] Ugander, J., Karrer, B., Backstrom, L., and Marlow, C., The anatomy of the Facebook social graph, preprint arxiv:1111.4503 (2011).
- [453] Valverde, S., Cancho, R. F., and Solé, R. V., Scale-free networks from optimal design, *Europhys. Lett.* **60**, 512–517 (2002).
- [454] van den Heuvel, M. P. and Hulshoff Pol, H. E., Exploring the brain network: A review on resting-state fMRI functional connectivity, *Eur. Neuro-psychopharmacol.* **20**, 519–534 (2010).
- [455] van Hooff, J. A. R. A. M. and Wensing, J. A. B., Dominance and its behavioral measures in a cap-

- tive wolf pack, in H. Frank, ed., *Man and Wolf*, pp. 219–252, Junk Publishers, Dordrecht (1987).
- [456] Vázquez, A., Flammini, A., Maritan, A., and Vespignani, A., Modeling of protein interaction networks, *Complexus* **1**, 38–44 (2003).
- [457] Vázquez, A. and Moreno, Y., Resilience to damage of graphs with degree correlations, *Phys. Rev. E* **67**, 015101 (2003).
- [458] Vázquez, A., Pastor-Satorras, R., and Vespignani, A., Large-scale topological and dynamical properties of the Internet, *Phys. Rev. E* **65**, 066130 (2002).
- [459] von Mering, C., Jensen, L. J., Snel, B., Hooper, S. D., Krupp, M., Foglierini, M., Jouffre, N., Huynen, M. A., and Bork, P., STRING: Known and predicted protein-protein associations, integrated and transferred across organisms, *Nucleic Acids Res.* **33**, D433–D437 (2005).
- [460] Wagner, C. S. and Leydesdorff, L., Network structure, self-organization, and the growth of international collaboration in science, *Res. Policy* **34**, 1608–1618 (2005).
- [461] Wang, D. J., Shi, X., McFarland, D. A., and Leskovec, J., Measurement error in network data: A reclassification, *Soc. Networks* **34**, 396–409 (2012).
- [462] Wasserman, S. and Faust, K., *Social Network Analysis*, Cambridge University Press, Cambridge (1994).
- [463] Watts, D. J., *Small Worlds*, Princeton University Press, Princeton (1999).
- [464] Watts, D. J., A simple model of global cascades on random networks, *Proc. Natl. Acad. Sci. USA* **99**, 5766–5771 (2002).
- [465] Watts, D. J., Dodds, P. S., and Newman, M. E. J., Identity and search in social networks, *Science* **296**, 1302–1305 (2002).
- [466] Watts, D. J. and Strogatz, S. H., Collective dynamics of ‘small-world’ networks, *Nature* **393**, 440–442 (1998).
- [467] West, D. B., *Introduction to Graph Theory*, Prentice Hall, Upper Saddle River, NJ (1996).
- [468] West, G. B., Brown, J. H., and Enquist, B. J., A general model for the origin of allometric scaling laws in biology, *Science* **276**, 122–126 (1997).
- [469] West, G. B., Brown, J. H., and Enquist, B. J., A general model for the structure and allometry of plant vascular systems, *Nature* **400**, 664–667 (1999).
- [470] White, J. G., Southgate, E., Thompson, J. N., and Brenner, S., The structure of the nervous system of the nematode *Caenorhabditis elegans*, *Phil. Trans. R. Soc. London B* **314**, 1–340 (1986).
- [471] Wilf, H., *Generatingfunctionology*, 2nd edn., Academic Press, London (1994).
- [472] Willis, J. C. and Yule, G. U., Some statistics of evolution and geographical distribution in plants and animals, and their significance, *Nature* **109**, 177–179 (1922).
- [473] Wodak, S. J., Pu, S., Vlasblom, J., and Séraphin, B., Challenges and rewards of interaction proteomics, *Mol. Cell. Proteomics* **8**, 3–18 (2009).
- [474] Wuellner, D. R., Roy, S., and D’Souza, R. M., Resilience and rewiring of the passenger airline networks in the United States, *Phys. Rev. E* **82**, 056101 (2010).
- [475] Yan, X., Bayesian model selection of stochastic block models, in *Proceedings of the 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pp. 323–328, Institute of Electrical and Electronics Engineers, New York (2016).
- [476] Yang, Z., Algesheimer, R., and Tessone, C. J., A comparative analysis of the community detection algorithms on artificial networks, *Sci. Rep.* **6**, 30750 (2016).
- [477] Yook, S. H., Jeong, H., and Barabási, A.-L., Modeling the Internet’s large-scale topology, *Proc. Natl. Acad. Sci. USA* **99**, 13382–13386 (2001).
- [478] Yule, G. U., A mathematical theory of evolution based on the conclusions of Dr. J. C. Willis, *Philos. Trans. R. Soc. London B* **213**, 21–87 (1925).
- [479] Zachary, W. W., An information flow model for conflict and fission in small groups, *J. Anthropol. Res.* **33**, 452–473 (1977).
- [480] Zanette, D. H. and Manrubia, S. C., Vertical transmission of culture and the distribution of family names, *Physica A* **295**, 1–8 (2001).
- [481] Zhang, P. and Moore, C., Scalable detection of statistically significant communities and hierarchies, using message passing for modularity, *Proc. Natl. Acad. Sci. USA* **111**, 18144–18149 (2014).
- [482] Zhang, X., Martin, T., and Newman, M. E. J., Identification of core-periphery structure in networks, *Phys. Rev. E* **91**, 032803 (2015).

REFERENCES

- [483] Zhang, X., Moore, C., and Newman, M. E. J., Random graph models for dynamic networks, *Eur. Phys. J. B* **90**, 200 (2017).
- [484] Zhang, X. and Newman, M. E. J., Multiway spectral community detection in networks, *Phys. Rev. E* **92**, 052808 (2015).
- [485] Zheng, X., Zeng, D., Li, H., and Wang, F., Analyzing open-source software systems as complex networks, *Physica A* **387**, 6190–6200 (2008).
- [486] Zipf, G. K., *Human Behaviour and the Principle of Least Effort*, Addison-Wesley, Reading, MA (1949).

INDEX

Page numbers in **bold** denote definitions or principal references.

- 1-component 133, **181**, 304
- 2-component 181
 - giant 367
 - random graph 367
- 3-component 181

- acquaintance immunization **592**, 605
- actors 47, 105
 - film actor network **61**, 115, 117, 171, 176
- acyclic network 38, **111**
 - absence of loops in 38, **111**
 - adjacency matrix 113
 - allowed degree sequence 420
 - approximately acyclic 38, 40, 41, 98, 111, 564
 - cascade model 368
 - citation network **38**, 40, 111
 - components 309
 - degree sequence 420
 - food web 98
 - models **419**, 438
 - Price model 438
 - random graph model 419
 - rank structure 564
 - strongly connected
 - components 309
 - test for 112
 - visualization 112
- Add Health study **54**, 201
- adjacency list 229
 - adding an edge to 232
 - computational complexity **232**, 234
 - directed network 230
 - enumerating neighbors 233
 - finding an edge in 232
 - hash tables 234
 - hybrid representation 233
 - multiedges 230
 - removing an edge from 232
 - self-edges 230
 - undirected network 229
- adjacency matrix **106**, 226
 - acyclic network 113
 - adding an edge to 226
 - asymmetric 110
 - computational complexity of operations **226**, 227
 - computer representation 226
 - diagonal elements **107**, 130
 - directed network 110
 - disadvantages 227
 - eigenvalues 160, 161, 164, 661, 690, 697, **698**, 708
 - eigenvectors 160, 161, 211, 647, 690
 - hybrid representation 233
 - interlayer 121
 - largest eigenvalue 164, 661, **698**, 700, 708
 - multiedges **107**, 110, 130
 - multilayer network 120
 - multiplex network 120
 - negative elements 109
 - powers of 131
 - projection of bipartite network 118
 - removing an edge from 226
 - self-edges **107**, 110
 - simple network **106**, 131
 - sparse network 227
 - spectrum 698
 - strictly triangular 113
 - symmetric **107**, 226
 - triangular 113
 - undirected network **106**, 226
 - weighted network **108**, 226
- affiliation network 49, **60**, 114
 - bipartite network **61**, 114
 - boards of directors 61
 - CEOs of companies 61
 - coauthorship 61
 - film actors 61
 - Southern Women Study **49**, 60
- affinity purification 80
 - errors in 280
- aggregated nodes 277
- airline network 3, **28**, 118, 479, 480
 - edge lengths 109, **563**
 - hub-and-spoke **479**, 482
 - models 480
 - optimization **479**, 480, 486
 - star graph 481
- algebraic connectivity 152
 - and visualization 147
- algorithm 218
 - augmenting path algorithm 263
 - average distance 249
 - average-linkage clustering 535
 - maximum flow 262
 - Barabási–Albert model 450
 - basic network quantities 237

- betweenness centrality **252**, 530
- breadth-first search **34**, **241**
- CFinder 551
- closeness centrality 249
- clustering coefficient 239
- community detection 498
- complete-linkage clustering 535
- computational complexity of **221**, **223**
- connectivity **262**, **269**
- cut set 252
- degree 238
- diameter 249
- Dijkstra's algorithm 258
- edge cut set **262**, **268**
- edge-independent paths **262**, **268**
- Edmonds–Karp algorithm **263**
- EM algorithm 286
- expectation–maximization algorithm 286
- Ford–Fulkerson algorithm **262**
- greedy algorithm 719
- heuristic **500**, 519
- hierarchical clustering 534
- HITS algorithm 168
- independent paths **262**, **268**, **269**
- InfoMap **515**, 550
- k*-cores 179
- Louvain algorithm 511
- maximum flow 262
- minimum cut **262**, **268**
- minimum violations ranking 564
- modularity maximization 498
- network visualization 145
- node cut sets 268
- node-independent paths 269
- percolation 594
- preferential attachment **443**, **450**
- preflow-push algorithm **263**
- Price model **443**
- proof of correctness 241
- running time 221, 223
- shortest distance **241**, **247**, **257**
- shortest path **241**, **247**, **249**, **257**
- single-linkage clustering 535
- site percolation 594
- spectral community detection 505
- time complexity **221**, **223**
- visualization 145
- allometric scaling 31
- alter 55
- animal networks 49, **58**, 279, 565
 - dominance hierarchies **58**, 565
 - errors in 279
 - rank structure in 565
- Antarctic food web 96
- arcs 110
- Arts and Humanities Citation Index 37
- assortative mixing **201**, 203, 206, 209, **335**
 - and core–periphery structure 210
 - and latent-space structure 559
- assortativity coefficient **209**, 211, 336
 - by age 206
 - by degree **209**, 336
 - by ethnicity 201
 - by gender 48, 202
 - by income 206
 - by location 209
 - by ordered characteristics 206
 - by unordered characteristics 203
 - by vector characteristics 209
- correlation coefficient **209**, 336
- disassortative mixing 202, 336
- friendship network **201**, 206, 209, 336
 - geographic 209
 - measures of **201**, 203, 206, 209
 - models of 420, 421
 - Pearson correlation **209**, 336
 - random graph models 420, 421
 - real-world networks 335
 - social networks 54, 55, **201**, 206, 335
 - stochastic block model 421
 - World Wide Web 203
- assortativity coefficient **209**, 211, 336
 - calculation of 336
- attracting fixed point **679**, 681, 688
- augmenting path 264
- augmenting path algorithm 262
 - computational complexity 266
 - implementation 265
 - proof of correctness 266
 - worst case performance 263
- authority centrality 168
- autonomous system 22
- average degree **127**, 130
 - bipartite networks 155
 - directed networks 130
 - friendship networks 129
 - giant component 431
 - Internet 380
 - neighbors 379
 - Poisson random graph 345
 - preferential attachment 456
 - Price model 456
 - random graph 345
 - small components 358, 392
 - tree 154
 - two-mode network 155
 - undirected network 127
- average distance 170, 249, **311**
 - algorithm for 249
 - all pairs 249, **311**
 - closeness centrality **170**, 249, 331
 - harmonic mean 172
 - log *n* behavior 312
- average-linkage clustering 535
- AVL tree 122, 234
- axon 88
- balance theorem 192
 - proof 216
- Barabási–Albert model 448
 - addition of extra edges 459
 - computer simulation 450
 - degree distribution 450
 - exponent 450

- generalizations 458
- master equation 449, 460
- non-linear attachment 466
- power-law distribution **450**, 455, 471
- relation to Price model 449
- removal of edges 461
- removal of nodes 464, **466**
- solution 449
- time evolution 458
- without growth 488
- without preferential attachment 490
- basic reproduction number 616
 - SIR model 617
 - SIS model 618
- Bernoulli random graph **343**, 345
- beta function 441
 - integral form 448
 - power-law tail 442
- Bethe lattice 122, 575
- betweenness centrality 173
 - algorithm for **252**, 530
 - community detection
 - algorithm **530**, 553
 - computational complexity 255, 257
 - directed networks 175
 - distribution 330
 - edge betweenness 530
 - film actor network 176
 - flow betweenness 177
 - Internet 330
 - normalization 176
 - power-law distribution 331
 - random-walk betweenness 177
 - star graph 215
 - tree 213
 - variants 177
- BGP **17**, 21
- BGP table 21
- bibliographic coupling 39
- bibliographic coupling network 39
- bibliometrics 37
- bicomponents 181
 - random graph 367
- bifurcations 692
- big-O notation 222
- binary heap **261**, 536
- binary tree 553, 725
 - dendrogram 531
 - hierarchical network model 724
- biochemical network **6**, **70**
- biological network **5**, **70**
 - animal social network 49, **58**, 279, 565
 - biochemical networks **6**, **70**
 - blood vessel network 31
 - brain connectivity 94
 - disease network 87
 - drug interaction network 86
 - drug-target network 87
 - duplication-divergence process 474
 - ecological networks **5**, **95**
 - errors in 279
 - food webs **5**, **96**
 - functional connectivity 94
 - genetic regulatory networks **6**, **80**
 - host-parasite network 100
 - metabolic networks **6**, **70**
 - mutualistic networks 100
 - neural networks **5**, **88**
 - node copying 474
 - protein-protein interaction networks **6**, **76**
 - protein structure networks 87
 - RNA structure network 87
 - root network 31
 - scale-free 477
 - seed dispersal networks 100
- bipartite network 115
 - affiliation network **61**, 114
 - average degree 155
 - boards of directors 61
 - CEOs of companies 61
 - coauthorship 61
 - directed 72, 116
 - disease network 87
 - drug-target network 87
 - film actors **61**, 115, 117
 - incidence matrix 116
- keyword index 45
- metabolic networks **72**, 116
- models for **419**, 433
- mutualistic networks 101
- one-mode projection 29, 72, **116**, 178, 419
- protein-protein interaction networks 77
- rail network 29
- random graph model **419**, 433
- recommender networks **44**, 115
- sexual contact networks 116
- social networks 61
- Southern Women Study **49**, 60
- weighted 116
- blogs 60
 - political 540
- blood vessel network 31
- boards of directors network 61
- bond 106
- bond percolation **571**, 627
 - and SIR model **627**, 629
 - configuration model 629
 - disease spreading 627
 - giant cluster 629
 - joint site/bond 672
 - percolation threshold 630
 - random graph 604
 - scale-free network 633
 - square lattice 603
- bootstrap percolation 179, 639
- Border Gateway Protocol **17**, 21
- bow tie diagram 309
- Brainbow 93
- brain cells 88
- brain connectivity network 94
- brain network 88
 - Brainbow 93
 - Clarity 93
 - empirical measurement 90, 94
 - functional connectivity 94
 - functional MRI 95
 - macroscopic 94
 - neurons 88
- breadth-first search **34**, **241**
 - augmenting path algorithm 263

INDEX

- betweenness centrality 252
 - closeness centrality 249
 - computational complexity 222, **246**
 - finding components 134
 - finding percolation clusters 594
 - finding shortest paths 249
 - implementation 243, **244**, 249
 - peer-to-peer network 714
 - proof of correctness **242**, 273
 - running time 222, **246**
 - sparse networks **247**, 249
 - two-source **247**, 252
 - using hash tables 249
 - variants 247
 - web crawling 34
- broker 176
- Caenorhabditis elegans*, see *C. elegans*
- call graph
 - software 25
 - telephone 59
- cascade model 368
- cascading failure **28**, 608
- cavity method 358
- Cayley tree 122, 338, 575
- C. elegans*
 - metabolic network 477
 - neural network 92
- cellphone data 59, 293
- centrality 9, **159**
 - authority centrality 168
 - betweenness centrality 173
 - closeness centrality 170
 - degree centrality 159
 - distribution 330
 - edge betweenness 530
 - eigenvector centrality **159**, 330, 648, 661, 668, 730
 - flow betweenness 177
 - HITS algorithm 168
 - hub centrality 168
 - hubs and authorities 168
 - in-degree 159
 - Katz centrality 163
 - out-degree 159
 - PageRank 165
 - random-walk betweenness 177
 - regular networks 211
- CEO network 61
- CFinder algorithm 551
- chromatic number 125
- Chung–Lu model 375, **377**, 423
- circuit-switched network 25
- citation counts 159
- citation network 4, **37**, 39
 - academic papers 37
 - acyclic **38**, 40, 111
 - Arts and Humanities Citation Index 37
 - bibliographic coupling 39
 - Citebase 38
 - CiteseerX 38
 - cocitation 39
 - cumulative degree distribution 322
 - data 37
 - degree distribution 322
 - errors in 37
 - in-degree distribution 322
 - legal citations 41
 - loops in **38**, 40, 111
 - models for **435**, 472
 - node copying model 472
 - node ordering 111
 - out-degree distribution 322
 - papers 37
 - patents 39
 - power-law degree distribution 37, 39, 317, 322, **435**, 443
 - preferential attachment model **435**, 443, 472
 - Price model **435**, 443
 - rank structure 564
 - scale-free network 37, 39, 317, 322, **435**, 443
 - Science Citation Index 37
 - Scopus 37
 - Social Science Citation Index 37
 - statistics 38
 - visualization 112
- Citebase 38
- CiteseerX 38
- clique 178
 - in one-mode projection **117**, 178
 - transitivity 183
- closed triad **183**, 333
- closeness centrality 170
 - algorithm for 249
 - computational complexity 249
 - distribution 331
 - film actor network 171
 - harmonic mean 172
 - Internet 331
 - tree 213
- cluster 192, 573
 - giant **572**, 574, 580, 587, 590, 599, 627, 629
 - percolating 573
 - spanning 573
- clusterability **192**, 216
- clustering
 - clustering coefficient **183**, 332
 - community detection **498**, 534
 - complete linkage 535
 - hierarchical 534
 - local **186**, 334
 - single linkage 535
- clustering coefficient **183**, 332
 - algorithm for 239
 - coauthorship network 185, 333
 - computational complexity 239
 - configuration model 332, **381**
 - directed network 185
 - email network 185
 - film actor network 185
 - food web 334
 - Internet 333
 - local clustering coefficient **186**, 334
 - network model 368, 421, 425
 - random graph 332, 334, **347**, 364, 366, 381
 - real-world networks 332
 - social networks **185**, 333
 - tree 184
 - Watts–Strogatz variant **188**, 334
 - World Wide Web 334

- coauthorship network 49, **61**, 495
 - affiliation network 61
 - average neighbor degree 380
 - bipartite network 61
 - clustering coefficient 185, 333
 - friendship paradox 380
 - funneling 313
 - percolation on **600**, 602
 - transitivity 185, 333
- cocitation 39
- coexistence threshold 636
- co-immunoprecipitation 77
 - errors in 280
- coinfection 622, **637**
 - configuration model 637
 - epidemic threshold 639
 - on a random graph 638
- co-links 189
- collaboration network 49, **61**, 495
 - affiliation network 61
 - average neighbor degree 380
 - bipartite network 61
 - clustering coefficient 185, 333
 - friendship paradox 380
 - funneling 313
 - percolation on **600**, 602
 - transitivity 185, 333
- collaborative filtering 44
- college football network 541
- coloring 125
- common neighbors 195
- community detection **494**, 498
 - algorithms 498
 - and network visualization 496
 - average-linkage clustering 535
 - benchmarks **538**, 539, 542, 543
 - best methods 550
 - betweenness algorithm **530**, 553
 - bisection 499
 - CFinder algorithm 551
 - clustering **498**, 534
 - comparison of algorithms 550
 - complete-linkage clustering 535
 - dendrogram 531
 - edge clustering 552
 - four groups test 543
 - fraction of nodes correct 544
 - heuristic algorithms **500**, 519
 - hierarchical clustering 553
 - information theory 515
 - Jaccard coefficient 197
 - LFR benchmark network 543
 - loop counting algorithm 532
 - Louvain algorithm 511
 - map equation 519
 - maximum likelihood 520
 - modularity matrix 505
 - modularity maximization 498
 - more than two communities 509
 - mutual information 547
 - normalized mutual information 547
 - number of communities 509, 529
 - overlapping communities 551
 - performance 538, 544
 - permutation of groups 545
 - profile likelihood 528
 - random walk algorithm **515**, 553
 - repeated bisection 509
 - single-linkage clustering **535**, 537
 - spectral algorithm **505**, 511
 - statistical inference 520, **522**, 551, 553
 - stochastic block model 423, **522**, 542
 - tests of 539
 - two communities 499
- community food web 99
- community structure 10, **494**
 - and visualization 496
 - dendrogram 531
 - detection **494**, 498
 - friendship networks 496, 503
 - in random graphs 365
 - karate club network **503**, 507, 537, **539**
 - metabolic network 496
 - models of **421**, 522
 - multilayer networks 552
 - mutual information 548
 - normalized mutual information **547**, 549
 - overlapping communities 551
 - permutation of groups 545
 - planted partition model 543
 - social networks 201, 206, 495, 503, 539
 - World Wide Web 10, 177, 496, 540
- compartmental model 609
- complementary cumulative distribution function 321
- complete-linkage clustering 535
- complex contagion **623**, 639
 - and *k*-cores 640
 - avalanches 645
 - epidemic threshold 644
 - initial conditions 640
 - outbreak size 642
 - phase transition 644
 - Poisson random graph 642
- components **133**, 304
 - algorithm for 243
 - average size 391
 - acyclic network 309
 - algorithm for finding 134
 - bicomponents **181**, 367
 - configuration model 384, 391, 411, 414
 - directed networks 36, **134**, 308
 - disease spreading 625
 - extensive **348**, 354
 - film actor network 306
 - giant 306, **348**, 384
 - in-components **136**, 309
 - Internet 307
 - k*-components 180
 - k*-connected 180
 - largest 306
 - out-components **136**, 308
 - random graphs **347**, 354, 355, 384, 391, 411, 414
 - real-world networks 304
 - small components **306**, 308
 - strongly connected **135**, 308

INDEX

- tricomponents 181
- undirected networks **133**, 180, 306
- weakly connected **135**, 308
- World Wide Web 36, 134, **308**
- compressed exponential 471
- computational complexity 221
 - adjacency list operations **232**, 234
 - adjacency matrix operations 226
 - augmenting path algorithm 266
 - betweenness centrality 255, 257
 - breadth-first search 222, **246**
 - closeness centrality 249
 - clustering coefficient 239
 - diameter 249
 - Dijkstra's algorithm 260
 - modularity maximization 504, 507, 512
 - on sparse networks **223**, 227, 232
 - percolation algorithms **594**, 598, 599
 - reciprocity 274
 - shortest distance **246**, 249, 260
 - shortest path 246, 249, 251, 260
- computer algorithm, *see* algorithm
- conditional entropy 547
- configuration model 369
 - average component size 391
 - average neighbor degree 379
 - bipartite **419**, 433
 - bond percolation on 629
 - clustering coefficient 332, **381**
 - components 384, 391, 411, 414
 - cross-immunity model 634
 - definition 370
 - diameter 399
 - directed **417**, 433
 - edge probability **373**, 376
 - epidemic threshold **630**, 666, 670
 - excess degree distribution **377**, 381
 - with given expected degrees 375
- connectance 128
- connected network 121, **134**
- connected triple 185
- connectivity 138
 - algebraic 152
 - algorithm for 263
 - edge connectivity **138**, 262
 - node connectivity **138**, 269
 - robustness **181**, 262
- consumer ISP 16
- contact tracing 67
- continuous phase transition 582
- continuous-time dynamical system 675
- core 178, 555
 - assortative mixing 210
 - core-periphery structure 180, 555
 - k*-core 178
- core-periphery structure 180, **555**
 - algorithms for 555
 - assortative mixing 210
 - Borgatti–Everett algorithm 556
 - detection 555
 - k*-cores **180**, 556
 - modularity function for 556
 - statistical inference 558
 - stochastic block model 558
- correlation coefficient
 - assortative mixing **209**, 336
 - degree correlation **211**, 238, 336
 - rows of adjacency matrix 197
 - similarity measure 197
- cosine similarity 196
 - and community detection 534, 537
 - and disambiguation 300
 - and entity resolution 300
 - and hierarchical clustering 534, 537
 - and link prediction 299
- coupled oscillators 701
- crawler, *see* web crawler
- cross-immunity 621
- excess degree 380
- exponential degree
 - distribution 416, 431, 580, 605
- friendship paradox 379
- geometric degree distribution 416, 431, 580, 605
- giant component **384**, 431
- giant percolation cluster **574**, 580, 582
- graphical solution **389**, 575
- k*-regular 429, 602, 671, 673
- largest eigenvalue 700
- locally tree-like 382
- mean-square component size 431
- multiedges **371**, 373
- neighbors at a given distance **384**, 411
- non-uniform percolation **587**, 604
- number of multiedges 373
- number of self-edges 374
- pair approximation 653
- percolation on **574**, 586, 602, 603, 604, **629**, 633, 637
- percolation threshold **577**, 586, 590, 602, 603, 604, **630**
- phase transition 385, 574, 577, 630
- Poisson degree distribution **372**, 407
- power-law degree distribution 372, **395**, 578, 585, 633
- regular network 429, 602, 671, 673
- scale-free network 372, **395**, 578, 585, 633
- second neighbors **383**, 408
- self-edges **371**, 374
- SI model on **655**, 673
- SIR model on **629**, 662, 671
- site percolation on **574**, 586, 602, 603, 604
- small components **391**, 411, 414
- small-world effect 401
- third neighbors **384**, 411

- coexistence threshold 636
- configuration model 634
- epidemic threshold 636
- random graph 637
- cross-validation 298
- cumulative advantage 435
- cumulative degree distribution 321
 - calculation of 322
 - citation network 322
 - complementary 321
 - Internet 322
 - power law **321**, 447
 - preferential attachment model 447
 - Price model 447
 - scale-free network **321**, 447
 - World Wide Web 322
- cut set 137, **139**
 - algorithm for 268, **269**
 - and robustness **181**, 262
 - edges **139**, 141, 262, **268**
 - minimum **139**, 262
 - nodes 137, **139**, 181
 - weighted network 141
- cycle 111
- DAG 111
- data packet 3, 10, **15**, 17, 21
- data structures 225
- dating network 49, 60
- Davis's balance theorem **193**, 216
- degree 9, **126**, 159
 - algorithm for 238
 - assortative mixing by **209**, 336
 - average **127**, 130
 - calculation 238
 - centrality measure 159
 - correlation coefficient 336
 - correlations 336
 - directed networks **130**, 316
 - disassortative mixing by 336
 - distribution **313**
 - excess degree 380
 - friendship networks **52**, 126, 129
 - homophily by **209**, 336
 - in-degree 9, **130**
 - in small components 358, 392
 - mean **127**, 130
 - mean square 240, 327, 328, 578
 - notation for **127**, 438
 - of neighbor 379
 - out-degree 9, **130**
 - second moment 240, 327, 328, 578
 - sequence **314**, 370
 - sparse networks 129
 - undirected networks 126
- degree centrality 159
- degree-corrected stochastic block model **422**, 522, 543
- degree correlation **209**, 336
 - algorithm for 238, **336**
 - correlation coefficient **211**, 238, 336
 - model of 420
 - non-social networks 336
 - random graphs 364
 - social networks 336
- degree distribution 313
 - Barabási–Albert model 450
 - citation network 37, 39, 317, **322**, 435
 - configuration model 371
 - cumulative 321
 - directed networks 316
 - excess degree **377**, 381
 - exponential 408, 416, 488, 580, 588, 599, 603, 605, 671
 - generating functions **386**, 407, 574
 - geometric 408, 416, 488, 580, 588, 599, 603, 605, 671
 - histogram **315**, 316, 317, 319, 321
 - in giant component 431
 - in/out-degree distribution **316**, 418, 433
 - in small components 366, 393
 - Internet 20, 280, **315**, 317, 321, 365
 - logarithmic binning 320
 - log–log plot **317**, 319
 - mean **127**, 130
 - mean square 327, 328, 578
 - metabolic network 477
 - moments 240, **327**, 328, 578
 - out-degree 316, 322
 - peer-to-peer network 43
 - plots of **315**, 316, 317, 321, 322, 365
 - Poisson **346**, 365, 407, 578, 604, 631, 637, 642
 - power law 20, 240, 312, **317**, 372, 395, 403, 435, 442, 476, 578, 633
 - preferential attachment **438**, 446, 450, 461, 465, 468
 - Price model **438**, 446, 455
 - probability generating function **386**, 407, 574
 - random graph **346**, 365
 - real-world networks 313
 - right-skewed **316**, 365
 - scale-free network 20, 240, **317**
 - second moment 327, 328, 578
 - stochastic block model 422
 - stretched exponential 465, **470**
 - tail **315**, 318, 326, 404
 - World Wide Web 317, **322**, 329
- degree sequence 314, 370
- degrees of separation 10, **63**, 718
- delivery route network 29
- dendrite 88
- dendrogram 122, **531**, 553, 724
 - community detection 531, 537
 - karate club network 537
 - hierarchical clustering 537
 - hierarchical network model 724
- dense network 128
- density 128
- deterministic dynamical system 675
- diameter 133
 - algorithm for 249
 - computational complexity 249
 - configuration model 399
 - log n behavior 312, 360, **362**, 401
 - power-law degree distribution 312
 - random graph **360**, 399

INDEX

- real-world networks 312
 - scale-free network 312
 - small-world effect 133, 312, 360, **363**
 - square lattice 337
 - diffusion MRI 95
 - diffusion tractography 95
 - digraph, *see* directed network
 - Dijkstra's algorithm 258
 - computational complexity 260
 - finding shortest paths 261
 - implementation 260
 - proof of correctness 259
 - technological uses 262
 - using a heap 261
 - directed acyclic graph 38, **111**
 - directed edges 4, **110**
 - directed network 4, **110**
 - acyclic 111
 - adjacency list 230
 - adjacency matrix 110
 - average degree 130
 - betweenness centrality 175
 - bipartite 72, 116
 - citation networks 4, **37**
 - clustering coefficient 185
 - components 36, **134**, 308
 - configuration model for **417**, 433
 - correlation of in- and out-degree 274, **316**, 433
 - degree **130**, 316
 - degree distribution 316
 - dynamical system on 707
 - edge probability 417
 - eigenvector centrality 161
 - errors in 277
 - food webs 6, **96**, 110
 - friendship networks 52
 - graph Laplacian 143
 - in-degree 130
 - in-degree distribution 316
 - in/out-degree distribution **316**, 418, 433
 - latent-space structure 563
 - loops 111
 - metabolic networks 73
 - models of **416**, 433
 - multiedges 110, 130
 - number of edges 130
 - out-degree 130
 - out-degree distribution 316
 - paths in 312
 - random graph models 367, **416**, 433
 - reciprocity 189
 - rank structure **558**, 564
 - self-edges **110**, 130
 - shortest paths 312
 - small components 308
 - social networks **52**, 185
 - transitivity 185
 - trees 121
 - World Wide Web 4, **33**, 35, 110
 - directors, boards of 61
 - disambiguation 277, **300**
 - disassortative mixing 202
 - by degree 336
 - by gender 202
 - modularity 205
 - stochastic block model 422
 - disconnected network 134
 - discontinuous phase transition 582
 - complex contagion 644
 - discrete-time dynamical system 675
 - disease network 87
 - disease spreading 607
 - and eigenvector centrality **648**, 661, 668
 - and network robustness **579**, 590
 - and small-world effect 311
 - as a dynamical system 645, **676**, 678, 686, 688
 - basic reproduction number 616
 - coexistence threshold 636
 - combinations of diseases **620**, 633, 637
 - compartmental models 609
 - contact tracing 67
 - dynamics 645
 - endemic disease **618**, 620, 668, 670
 - equivalence to bond percolation 627
 - fully mixed approximation 609
 - herd immunity **571**, 579
 - immunization **570**, 586, 672
 - infected state 609
 - infection rate 610, 624
 - initial conditions **615**, 618, 640, 647, 660, 667, 669
 - logistic growth 611, 618
 - mass-action approximation 609
 - models of 608
 - MSIR model 620
 - on regular networks 671, 673
 - on scale-free networks 633
 - pair approximation 650
 - percolation theory **625**, 627
 - recovered state 612
 - recovery rate 613
 - reinfection **617**, 619
 - removed state 612
 - Reed–Frost model 672
 - SEIR model 620
 - SI model **609**, 625
 - SIR model **612**, 626
 - SIRS model 619
 - SIS model **617**, 667
 - social contagion 607, **622**, 639
 - susceptible state 609
 - transmission probability 626
 - vaccination **570**, 586, 672
- disjoint paths, *see* independent paths
- distributed databases 713
- distribution networks 29
 - optimization 479
- doctor network 53
- dolphin social network 58, **539**
- domains 21
- dominance hierarchy 58
 - rank structure 565
- drug interaction network 86
- drug–target network 87
- duplicate nodes 277
- duplication–divergence process 478
- dynamical system **675**, 685

- and graph Laplacian **691**, 697, 703
- bifurcations 692
- continuous time 675
- coupled oscillators 701
- deterministic 675
- discrete time 675
- disease spreading **645**, **676**, 678, 686, 688
- fixed points **677**, 686
- Floquet theory 706
- gossip model 693
- Jacobian matrix 680
- limit cycle **685**, 701
- linearization **677**, 686
- linear stability analysis **677**, 679, 686
- Lotka–Volterra equations 685
- multivariable 694
- non-network 676
- on directed networks 707
- on regular networks 706
- oscillating 683, 685, **701**
- oscillator networks 701
- SI model **646**, 676, 678, 685, 688
- SIR model 660
- SIS model 667
- stochastic 675
- synchronization 701
- dynamic network 118, **120**
 - multilayer network 120
 - Price model 451
 - random graph model 424
 - social networks 60, 120
- dynamic web page 35
- E. coli* 477
- ecological network 5, **95**
 - databases 100, 101
 - food web 5, **96**
 - host–parasite network 100
 - mutualistic network 100
 - plant networks 31, 100
 - plant–pollinator network 100
 - root network 31
 - seed dispersal 100
- edge betweenness centrality 530
- edge clustering 552
- edge connectivity 138
 - algorithm for 262
- edge cut set 139
 - algorithm for 262, **268**
 - weighted network 141
- edge-disjoint paths, *see* edge-independent paths
- edge-independent paths **137**, 262
 - algorithm for 262, **268**
 - and k -components 181
 - and maximum flow 140
 - and minimum cut **139**, 268
- edge list 236
- edge percolation, *see* bond percolation
- edge probability
 - Chung–Lu model 376
 - configuration model 373
 - directed configuration model 417
- edges 1, **106**, 108, 110
 - directed 110
 - errors on 277
 - failure of 571
 - hyperedges 114
 - lengths **109**, 257, 563
 - missing 277
 - multiedges **106**, 109, 110
 - negative 84, 89, 109, **190**, 258
 - number of **106**, 127, 130
 - percolation on 571, 627
 - reciprocated 189
 - removal of 461, 571
 - self-edges **106**, 107, 110
 - signed 84, 89, 109, **190**, 258
 - types of 53, **118**
 - weighted 108
- Edmonds–Karp algorithm 263
- ego 55
- ego-centered network 55
- eigenvalues
 - adjacency matrix 160, 164, 661, 690, 697, **698**, 708
 - graph Laplacian **150**, 701
 - largest 164, 661, **698**, 700, 708
 - Perron–Frobenius theorem **160**, 166, 169, 699
- Poisson random graph 708
- random graph 708
- regular network 708
- smallest 700
- square lattice 708
- star graph 154, 708
- eigenvector centrality **159**, 330
 - and infection probability **648**, 661, 668
 - and snowball sampling 66
 - and web crawling 730
- directed networks 161
- distribution 330
- Internet 330
- power-law distribution 330
- problems with 161
- regular network 211
- eigenvectors
 - adjacency matrix 160, 211, 647, 690
 - eigenvector centrality **159**, 330
 - graph Laplacian 151
 - Fiedler vector 561
 - leading **160**, 507, 648, 661, 668
 - modularity matrix 507
 - Perron–Frobenius theorem **160**, 166, 169, 699
 - square lattice 708
- electrical network
 - power grid 27
 - circuits 25, 148
 - resistor networks 148
- electricity grid, *see* power grid
- email network **59**, 190
 - clustering coefficient 185
 - errors in 279
 - reciprocity 190
 - small-world experiment on 64
 - transitivity 185
- EM algorithm 286
- embedding 559
 - and graph Laplacian 559
 - and network visualization **145**, 562
 - quality function for 559
 - statistical inference 563
- endemic disease **618**, 620, 668, 670

INDEX

- entity resolution 277, 300
- enzymes 71
- epidemic threshold 616, 628
 - and leading eigenvalue 661
 - and percolation threshold 627
 - complex contagion 644
 - configuration model 630, 666, 670
 - cross-immunity 636
 - Poisson degree distribution 631
 - power-law degree distribution 633
 - random graph 631
 - scale-free network 633
 - SIR model 616, 617, 630, 632, 661, 666
 - SIS model 618, 668, 670
- epidemics 607
 - and eigenvector centrality 648, 661, 668
 - and network robustness 579, 590
 - and small-world effect 311
 - as a dynamical system 645, 676, 678, 686, 688
 - basic reproduction number 616
 - coexistence threshold 636
 - combinations of diseases 620, 633, 637
 - compartmental models 609
 - contact tracing 67
 - dynamics 645
 - endemic disease 618, 620, 668, 670
 - equivalence to bond
 - percolation 627
 - herd immunity 571, 579
 - fully mixed approximation 609
 - immunization 570, 586, 672
 - infected state 609
 - infection rate 610, 624
 - initial conditions 615, 618, 640, 647, 660, 667, 669
 - logistic growth 611, 618
 - mass-action approximation 609
 - models of 608
 - MSIR model 620
 - on regular networks 671, 673
 - on scale-free networks 633
 - pair approximation 650
 - percolation theory 625, 627
 - recovered state 612
 - recovery rate 613
 - reinfection 617, 619
 - removed state 613
 - Reed–Frost model 672
 - SEIR model 620
 - SI model 609, 625
 - SIR model 612, 626
 - SIRS model 619
 - SIS model 617, 667
 - social contagion 607, 622, 623, 639
 - susceptible state 609
 - transmission probability 626
 - vaccination 570, 586, 672
- Erdős, Paul 345, 350
- Erdős–Rényi model, *see* random graph
- error correction 297
- errors 275
 - animal networks 279
 - biological networks 279
 - citation networks 37
 - contact tracing 67
 - directed networks 277
 - edge errors 277
 - email networks 279
 - error correction 297
 - estimation of 281, 285
 - exponents 325
 - false discovery rate 302
 - genetic regulatory networks 279
 - impact 276
 - Internet 18, 280
 - maximum likelihood method 282, 286
 - metabolic networks 279
 - missing data 277
 - models for 285, 290, 296
 - node errors 276
 - non-network data 281
 - protein–protein interaction
 - networks 79, 279
 - random-walk sampling 68
 - respondent-driven sampling 69
 - small-world experiment 63, 64
 - snowball sampling 66
 - social networks 52, 54, 63, 66, 67, 276, 278, 279, 293
 - sources of 278
 - surveys 54, 66, 278
 - tandem affinity purification 280
 - true positive rate 290, 294
 - two-hybrid screen 79, 280
 - types of 276
 - weighted networks 277
 - World Wide Web 35, 280
 - yeast two-hybrid screen 79, 280
- Escherichia coli* 477
- ethnographic studies 57
- Euclidean distance 197
- Euler characteristic 156
- excess degree 380
- excess degree distribution 377, 381
 - generating function 387, 407, 575
- expansion of a network 126
- expectation–maximization algorithm 286
- exponent 318
 - Barabási–Albert model 450
 - error on 325
 - formula for 324
 - Hill estimator 324
 - Internet 578
 - Lyapunov exponent 679
 - measurement 324
 - node copying model 476
 - power law 318, 324, 329, 442, 450, 476, 578
 - preferential attachment model 442, 450, 461
 - Price model 442

- statistical error on 325
 - World Wide Web 329
- exponential degree distribution
 - configuration model 416, 431, 580, 605
 - generating function 408
 - non-uniform percolation 588
 - percolation **580**, 588, 603, 605, 671
 - percolation threshold **581**, 603
 - power grid 599
 - small component sizes 416
- exponential distribution 403
 - generating function **403**, 408
 - stretched 465, 470
- exponential generating function 402
- extremely sparse network **129**, 223
- Facebook network 5, 49, **60**
 - small-world effect **63**, 312, 363
- face-to-face interaction network **57**, 59, 293
- faculty hiring network 565
- false discovery rate 302
- false positives and negatives 79, 277, 290, 298
- Fibonacci heap 261
- Fiedler vector 561
- file sharing network, *see*
 - peer-to-peer network
- film actor network 61
 - affiliation network 61
 - betweenness centrality 176
 - bipartite **61**, 115, 117
 - closeness centrality 171
 - clustering coefficient 185
 - largest component 306
 - most central actor 171, 176
 - one-mode projection 117
 - small-world effect 61
 - transitivity 185
- finite-size effects 584
- first-in/first-out buffer 245
- first mover effect 451, **456**, 487
- first-order phase transition **582**, 644
- fixed choice studies 52
 - problems with 279
- fixed point **677**, 686
 - attracting **679**, 681, 688
 - expansion around **677**, 686
 - mixed 679
 - neutral 679
 - non-symmetric 694
 - repelling **679**, 681, 688
 - saddle point **681**, 688
 - SI model 678, 688
 - symmetric **688**, 693
- Floquet theory 706
- flow betweenness 177
- fMRI 95
- food chain 96
- food web 5, **96**
 - acyclic 98
 - Antarctica 96
 - carbon flow 97
 - cascade model 368
 - clustering coefficient 334
 - community food web 99
 - databases 100
 - density 129
 - directed network 96, **97**, 110
 - edges 5, **96**
 - empirical measurement 99
 - energy flow **97**, 108
 - Little Rock Lake 5
 - model of 368
 - nodes 5, **96**
 - paleontological 99
 - rank structure 564
 - sink food web 99
 - source food web 99
 - transitivity 334
 - trophic level **98**, 156, 564
 - trophic species 96
 - weighted **100**, 108
- Ford–Fulkerson algorithm 262
- forest 121
- four-color theorem 124
- free choice studies **52**, 56, 279
- friendship network 5, 48, **51**
 - Add Health study **54**, 201
 - animosity 109
 - assortative mixing **201**, 206, 209, 336
 - average degree 129
 - community structure 496, 503
 - directed 52
 - edges 4, **47**, 106
 - Facebook 5, **60**, 363
 - homophily 201
 - in-degree 52
 - karate club network **5**, 57, **503**
 - node degrees **52**, 126, 129
 - nodes 4, **47**, 106
 - out-degree 52
 - rank structure 564
 - reciprocity 189
 - schoolchildren 48, 51, 53, 201, 206, 496
 - sparse 129
 - students **5**, 57, 293, **503**
- friendship paradox 379
 - coauthorship network 380
 - configuration model 379
 - generalized 429
 - Internet 380
- fully mixed approximation 609
 - equivalence to random graph 633
 - SI model **611**, 626, 676
 - SIR model **613**, 624, 626, 629, 632
 - SIRS model 620
 - SIS model 617
- functional connectivity network 94
- functional MRI 95
- funneling **64**, 313
 - coauthorship network 313
 - email experiments 64
 - Internet 313
 - real-world networks 313
- gas pipeline network 29
- gender
 - assortative mixing by 48
 - disassortative mixing by 202
- generalized friendship paradox 429
- General Social Survey 55
- generating function 386, 389, **401**
 - average of distribution 405
 - component sizes **411**, 414
 - definition 386, **402**

INDEX

- degree distribution **386**, 407, 574
- derivatives 402, **405**
- examples 402
- excess degree distribution **387**, 407, 575
- exponential distribution **403**, 408
- exponential generating function 402
- Fibonacci numbers 430
- first moment 405
- geometric distribution **403**, 408
- moments of distribution 405
- normalization 405
- Poisson distribution **403**, 407
- power-law distribution 397, **403**, 408
- powers of 406
- products of 406
- properties 404
- scale-free networks 397, **403**, 408
- small components **411**, 414
- random graph 407
- genes 83
 - duplication 478
 - expression 83
 - regulatory networks 80
- genetic regulatory network 6, **80**
 - databases 86
 - errors in 279
 - measurement of 84
- geodesic distance, *see* shortest distance
- geodesic path, *see* shortest path
- geographic structure 24
 - Internet 24
 - network of regions on a map 124
 - power grid 28
 - river networks 30
 - social networks 24, 59, 209, 559, 724
 - telephone network 27
 - transportation networks 28
- geometric degree distribution
 - configuration model 416, 431, 580, 605
 - generating function 408
 - non-uniform percolation 588
 - percolation **580**, 588, 603, 605, 671
 - percolation threshold **581**, 603
 - power grid 599
 - small component sizes 416
- geometric distribution 403
 - generating function **403**, 408
- Gephi (software package) 220
- giant component 306, **347**
 - average degree in 431
 - condition for 351, 384
 - configuration model **384**, 431
 - degree distribution in 431
 - graphical solution for 351, 389
 - lack of 308
 - more than one **307**, 354
 - power-law degree distribution 396
 - random graph **347**, 354, 616
 - real-world networks 306
 - scale-free network 396
 - size 306, **349**, 385
- giant in-component 308
- giant out-component 308
 - World Wide Web 36, 309
- giant percolation cluster 572
 - and epidemics 627
 - bond percolation **627**, 629
 - configuration model **574**, 580, 582, 587
 - exponential degree distribution 580
 - graphical solution for 575
 - real-world networks 599
 - scale-free networks **584**, 590, 600, 602, 633
 - site percolation 572
 - size **574**, 580, 629
- $G(n, m)$ 343
- $G(n, p)$ 344
- Gnutella 43
- Google 35, 710
 - PageRank 166, 712
- Google Patents 40
- Google Scholar 37
- gossip 10, **607**
 - models of 622, 693
- graph 1, **105**
- graph Laplacian 142
 - algebraic connectivity 152
 - and cut size 144
 - and dynamical systems **691**, 697, 703
 - directed networks 143
 - eigenvalues **150**, 701
 - eigenvectors 151
 - embedding 559
 - Fiedler vector 561
 - graph partitioning 143
 - largest eigenvalue 701
 - latent-space structure 560
 - multigraphs 143
 - network visualization 145
 - properties 150
 - random walks 147
 - resistor networks 148
 - sparsification 150
 - spectral gap 152
 - spectral partitioning 143
 - spectrum **150**, 701
 - weighted networks 143
 - zero eigenvalues **151**, 153
- graph partitioning 143
- graph theory 105
- Graphviz (software package) 220
- greedy algorithm
 - community detection **502**, 511
 - modularity maximization **502**, 511
 - network optimization 483
 - message passing 719
- Hamming distance 197
- Harary balance theorem 192
 - proof 216
- hash table 234
 - and adjacency list 234
 - and breadth-first search 249
 - load factor 235
- heap 261
 - binary heap **261**, 536

- Dijkstra's algorithm 261
- Fibonacci heap 261
- hierarchical clustering 536
- herd immunity 571, 579
- heuristic algorithm 500
- hidden population 65
- hierarchical clustering 534, 553
 - average-linkage clustering 535
 - complete-linkage clustering 535
 - karate club network 537
 - problems with 537
 - running time 536
 - single-linkage clustering 535, 537
 - using cosine similarity 534
- hierarchical random graph 553
- hierarchical structure 532, 534, 553, 564, 724
 - dendrogram 553
 - detection 530, 553, 564
 - directed networks 564
 - dominance hierarchy 58, 565
 - hierarchical clustering 534
 - hierarchical random graph 553
 - minimum violations ranking 564
 - network search 724
 - statistical inference 553, 564
- hierarchy 58, 564
- Hill estimator 324
- HITS algorithm 168
- homophily 201, 203, 206, 209, 335
 - assortativity coefficient 209, 211, 336
 - by age 206
 - by degree 209, 336
 - by income 206
 - by location 209
 - by ordered characteristics 206
 - by unordered characteristics 203
 - by vector characteristics 209
- correlation coefficient 209, 211, 238, 336
- friendship networks 201, 206, 209, 336
- geographic 209
- latent-space structure 559
- models of 421
- real-world networks 335
- social networks 54, 55, 201, 206, 335
- stochastic block model 421
- host-parasite network 100
- HTML (Hypertext Markup Language) 33
- HTTP (Hypertext Transfer Protocol) 33
- hub 9, 168, 315
 - airline networks 479, 482
 - hub centrality 168
 - hubs and authorities 168
 - random graphs 365
 - real-world networks 315
 - removal 272, 591
 - social networks 52
- hub-and-spoke network 479, 482, 486
- hub centrality 168
 - HITS algorithm 168
- hubs and authorities 168
- hyperedge 114
- hypergraph 114
- hyperlinks 3, 33, 105, 135, 280
- Hypertext Markup Language 33
- Hypertext Transport Protocol 33
- igraph (software library) 220
- immunization 570
 - acquaintance immunization 592, 605
 - and network robustness 579
 - by degree 572, 586, 590, 592
 - herd immunity 571, 579
 - network robustness 579
 - non-uniform 572, 586, 590, 592, 605
 - percolation theory 570, 572, 579, 586, 592, 605, 672
 - scale-free networks 579, 591
 - targeted 572, 586, 590, 592, 605
- incidence matrix 116
- in-components 136, 308
 - World Wide Web 281, 309
- in-degree 9, 130
 - average 130
 - centrality 159
 - citation network 322
 - correlation with out-degree 274, 316, 433
 - distribution 316, 317, 322, 329, 438
 - social networks 52
 - World Wide Web 9, 316, 317, 322, 329
- in-degree distribution 316
 - citation network 322
 - Price model 438
 - World Wide Web 316, 317, 322, 329
- independent paths 137, 262, 268, 269
 - algorithm for 262, 268, 269
 - and maximum flow 140
 - and minimum cut 139, 268
 - and robustness 181, 262
 - edge-independent 137, 262
 - Menger's theorem 139
 - node-independent 137, 269
 - number of 138
 - self-avoiding 131
 - uniqueness 138, 268
- INDEX experiments 65, 723
- infected state 609
- infection rate 610, 624
- infective state 609
- InFlow (software package) 220
- InfoMap algorithm 515, 550
 - running time 520
- information network 3, 32
 - citation networks 37, 39
 - keyword indexes 45
 - peer-to-peer networks 42, 713
 - recommender networks 44, 115
 - World Wide Web 3, 32
- information science 37
- information spreading 10, 607, 622, 693
- information theory 515
 - community detection 515
 - conditional entropy 547

INDEX

- mutual information 547
- normalized mutual information 547
- in/out-degree distribution **316**, 418, 433
- instant messaging network 59
- intermarriage network 58
- Internet 1, **15**
 - autonomous systems 20, **22**
 - average degree 380
 - average neighbor degree 380
 - backbone 15
 - betweenness centrality 330
 - closeness centrality 331
 - clustering coefficient 333
 - cumulative degree distribution 322
 - degree distribution 20, 280, **315**, 317, 321, 322, 365, 578
 - degree moments 328
 - domains 21
 - edges 15
 - eigenvector centrality 330
 - errors in data 280
 - exponent 578
 - failure of routers 569
 - friendship paradox 380
 - funneling 313
 - geographic structure 24
 - Internet Protocol 15
 - Internet service providers 15
 - largest component 307
 - measurement of 17, 280
 - network backbone providers 15
 - nodes **15**, 19, 22
 - packets **15**, 17
 - percolation on **569**, 571, 573, 578, 578, 600, 602
 - power-law degree distribution 20, 280, **317**, 321, 322, 578
 - robustness 181, **569**, 571, 573, 578, 592, 600, 602
 - router representation 19
 - routers **15**, 21, 569
 - scale-free network 20, 280, **317**, 321, 322, 578
 - small-world effect 10, 311
 - sparse network 129
 - subnets 20
 - susceptibility to attack 592
 - transitivity 333
 - weighted edges 108
- Internet Movie Database 171
- Internet packets **15**, 17
- Internet Protocol 15
- Internet routers 15
 - failure of 181, **569**, 571
 - routing tables 21
- Internet service providers 15
- interstate highway network 561
- interviews 51
 - problems with 54, 278
- IP address 15
 - and geographic location 24
 - subnets 20
- Ising spins 144
- ISP (Internet Service Provider) 15
 - consumer 16
 - regional 16
- Jaccard coefficient **197**, 299
- Jacobian matrix 680
- Jensen's inequality 287
- JUNG (software library) 220
- Karate Club Club 539
- karate club network 5, 57, **539**
 - community structure **503**, 507, 537, **539**
 - dendrogram 537
 - hierarchical clustering 537
 - modularity maximization **503**, 507
 - pictures of 5, 503
- Katz centrality 163
 - limiting values 212
 - parameter value 163
 - problems with 165
 - regular network 211
- Katz similarity 200
- k -club 182
- k -components 180
 - and node cut sets 181
 - and node-independent paths 180
 - and robustness 181
 - bicomponents **181**, 367
 - distinct from k -core 181
 - non-contiguous 182
 - random graph 367
 - tricomponents 181
- k -connected components, *see* k -components
- k -core 178
 - algorithm for 179
 - and bootstrap percolation 179
 - and complex contagion 640
 - distinct from k -component 181
- keyword index 45
- k -plex 182
- k -regular network, *see* regular network
- Kuratowski's theorem **126**, 155
- Lambert W -function 350, 363, 631
- Laplacian, *see* graph Laplacian
- latent Dirichlet allocation 46
- latent semantic analysis 46
- latent-space structure 558
 - and graph Laplacian 560
 - and visualization 562
- directed networks 563
- rank structure 564
- road network 561
- spectral algorithm 559
- statistical inference 563
- lattice
 - Bethe lattice 122
 - diameter 337
 - eigenvalues and eigenvectors 708
 - percolation on 603
 - regular network 128
 - square 337, 603, 708
- Lee, Christopher **171**, 176
- legal citation network 41
- letter-passing experiment, *see* small-world experiment
- LexisNexis 41
- LFR benchmark networks 543
- likelihood 283, 523

- community detection 520
- error estimation 282, 286
- normal distribution 282
- log-likelihood 284, 287, 521, 524
- maximum 282, 520
- profile likelihood 528, 567
- stochastic block model 523, 526
- limit cycle 685, 701
- linear stability analysis 677, 686
- LinkedIn 60
- link prediction 298
 - methods 298
 - using statistical inference 299, 554
 - using cosine similarity 299
- Little Rock Lake food web 5
- local clustering coefficient 186
 - and redundancy 188
 - average of 188, 334, 335
 - degree dependence 186, 335
 - real-world networks 334
- local ISP 16
- locally tree-like network 122, 353, 382, 383, 386, 391, 654
- logarithmic binning 320
- logistic equation 611
- logistic growth 611
 - SI model 611
 - SIS model 618
- log-likelihood 284, 287, 521, 524
- longitudinal network studies 60
- loops 132
 - acyclic networks 38, 111
 - and community detection 532
 - citation networks 38, 40, 111
 - counting 132
 - directed networks 111, 189
 - length two 189
 - length three 183, 239, 332
 - length four and longer 64
 - number of 132
 - random graphs 366
 - self-edges 106, 107, 110
 - signed networks 190
 - trees 121
 - triangles 183, 239, 332
- Lorenz curve 328, 338
- Lotka–Volterra equations 685
- Louvain algorithm 511
 - running time 512
- Lyapunov exponent 679
- map equation 519
- marriage network 58
- mass-action approximation 609
 - equivalence to random graph 633
 - SI model 611, 626, 676
 - SIR model 613, 624, 626, 629, 632
 - SIRS model 620
 - SIS model 617
- master equation 439
 - addition of extra edges 460
 - Barabási–Albert model 449
 - node copying model 475
 - non-linear attachment 467
 - preferential attachment model 439, 449, 451, 460, 462, 467
 - Price model 439, 451
 - removal of edges 462
 - without network growth 488
 - without preferential attachment 488
- master stability condition 691, 692
- master stability function 697, 704
- matrix
 - adjacency matrix 106, 226
 - graph Laplacian 142, 560, 691
 - incidence matrix 116
 - Jacobian 680
 - modularity matrix 505
- max-flow/min-cut theorem 140, 141
 - and random-field Ising model 142
 - on weighted networks 141
- maximum flow 140
 - algorithm for 262
 - and independent paths 140
 - and minimum cut 140, 262
 - augmenting path algorithm 263
 - max-flow/min-cut theorem 140, 141
- preflow-push algorithm 263
- weighted networks 141
- maximum likelihood
 - community detection 520
 - error estimation 282, 286
 - normal distribution 282, 284
- mean degree 127, 130
 - bipartite network 155
 - directed network 130
 - friendship network 129
 - giant component 431
 - Internet 380
 - neighboring node 379
 - Poisson random graph 345
 - preferential attachment 456
 - Price model 456
 - random graph 345
 - small components 358, 392
 - tree 154
 - two-mode network 155
 - undirected network 127
- measurement error 275
 - affinity purification 280
 - animal networks 279
 - biological networks 279
 - citation networks 37
 - contact tracing 67
 - directed networks 277
 - edge errors 277
 - email networks 279
 - error correction 297
 - estimation 281, 285
 - exponent values 325
 - false discovery rate 302
 - false positives and negatives 79, 277, 290, 298
 - genetic regulatory networks 279
 - impact 276
 - Internet data 18, 280
 - maximum likelihood methods 282, 286
 - metabolic networks 279
 - missing data 277
 - models for 285, 290, 296
 - node errors 276
 - non-network data 281

INDEX

- protein–protein interaction networks **79**, 279
- random-walk sampling **68**
- respondent-driven sampling **69**
- small-world experiment **63**, **64**
- snowball sampling **66**
- social networks **52**, **54**, **63**, **66**, **67**, **276**, **278**, **279**, **293**
- sources of **278**
- surveys **54**, **66**, **278**
- tandem affinity purification **280**
- true positive rate **290**, **294**
- two-hybrid screen **79**, **280**
- types of **276**
- weighted networks **277**
- World Wide Web **35**, **280**
- yeast two-hybrid screen **79**, **280**
- medical doctor network **53**
- Menger’s theorem **139**
- metabolic network **6**, **70**
 - bipartite **72**, **116**
 - C. elegans* **477**
 - community structure **496**
 - databases **76**
 - degree distribution **477**
 - directed **73**
 - E. coli* **477**
 - edges **6**, **71**
 - errors in **279**
 - measurement **74**
 - node copying **477**
 - nodes **6**, **71**
 - one-mode projection **72**
 - products **71**
 - scale-free network **477**
 - substrates **71**
- metabolic pathway **71**, **74**
- metabolic reaction **71**
- metabolite **71**
- Milgram experiment, *see* small-world experiment
- Milgram, Stanley **62**, **310**, **313**, **718**
- minimum cut **139**, **262**
 - algorithm for **262**, **268**
 - augmenting path algorithm **263**
 - and independent paths **139**, **268**
 - and maximum flow **140**, **262**
 - max-flow / min-cut theorem **140**, **141**
 - Menger’s theorem **139**
 - preflow-push algorithm **263**
 - weighted networks **141**
- minimum cut set **139**, **262**
 - algorithms for **268**, **269**
 - and *k*-components **181**
 - and robustness **181**, **262**
- edges **139**, **141**, **262**, **268**
- nodes **137**, **139**, **181**
- weighted networks **141**
- minimum spanning tree **122**
- minimum violations ranking **564**
- missing data **277**
 - missing edges **277**
 - missing nodes **276**
 - social networks **52**, **54**, **63**, **278**
- mixed fixed point **679**
- mixed-membership stochastic block model **552**
- mobile phone data **59**, **293**
- models
 - airline network **480**
 - acyclic networks **419**, **438**
 - assortative mixing **420**, **421**
 - Barabási–Albert model **448**
 - bipartite networks **419**, **433**
 - cascade model **368**
 - Chung–Lu model **375**, **377**, **423**
 - citation networks **435**
 - clustering **368**, **421**, **425**
 - community structure **421**, **522**
 - compartmental models **609**
 - configuration model **370**
 - degree-corrected stochastic block model **422**, **522**
 - degree correlations **420**
 - directed networks **416**, **433**
 - disease spreading **608**
 - duplication–divergence process **478**
 - dynamic networks **424**
 - epidemics **608**
 - Erdős–Rényi model **345**
 - food webs **368**
 - gossip **622**, **693**
 - growing networks **434**, **472**
 - hierarchical networks **553**, **724**
 - MSIR model **620**
 - network formation **434**
 - network optimization **479**
 - node copying **472**, **478**
 - planted partition model **543**
 - Poisson random graph **343**, **345**, **347**
 - preferential attachment **435**, **448**, **458**
 - Price model **435**
 - protein–protein interaction networks **477**
 - road network **486**
 - random graph **342**, **369**, **416**
 - Reed–Frost model **672**
 - SEIR model **620**
 - SI model **609**, **625**
 - SIR model **612**, **626**
 - SIRS model **619**
 - SIS model **617**, **667**
 - small-world effect **312**, **360**, **363**, **401**, **425**
 - small-world model **425**
 - stochastic block model **421**, **522**, **542**
 - transitivity **368**, **421**, **425**
 - voter model **622**
 - Watts–Strogatz model **425**
 - World Wide Web **435**, **438**, **458**, **459**, **466**
- modularity **204**, **498**
 - alternative forms **204**, **205**, **500**
 - community detection **498**
 - couples **216**
 - matrix form **505**
 - maximization **498**
 - negative **205**
 - sexual contact networks **205**
 - social networks **205**
- modularity matrix **505**

- community detection 505
- eigenvalues and eigenvectors 507
- modularity maximization 498
 - algorithms for 502, 505, 509, 511
 - comparison with other methods 550
 - computational complexity 504, 507, 512
 - extremal optimization 502
 - genetic algorithms 502
 - greedy algorithm 502, 511
 - karate club network 503, 507
 - Louvain algorithm 511
 - node moving algorithm 502, 511
 - resolution limit 512
 - simulated annealing 502
 - spectral algorithm 505, 511
- Molloy–Reed criterion 385
- moment closure 651
- Moreno, Jacob 47
- motifs 334
- movie actor network, *see* film actor network
- MRI 94
- MSIR model 620
- multiedges 106
 - configuration model 371, 373
 - directed networks 110, 130
 - in adjacency list 230
 - in adjacency matrix 107, 110, 130
 - preferential attachment 437, 439
 - Price model 437, 439
 - scale-free networks 374
 - self-edges 108
- multigraphs 106
 - and weighted networks 109, 141
- multilayer network 118
 - adjacency matrix 120
 - community structure in 552
 - dynamic networks 120
 - interlayer adjacency matrix 121
 - social networks 54, 120
 - transportation networks 118, 121
- multiplex network 119
 - adjacency matrix 120
 - social networks 54, 120
 - tensor representation 120
- mutual information 547
- mutualistic network 100
- name generator 51, 201
- Napster 42
- National Longitudinal Study of Adolescent Health 54, 201
- navigation 64, 718
- negative edges 84, 89, 109, 190, 258
- neighbors
 - at given distance 384, 411
 - average degree 379
 - common neighbors 195
 - second neighbors 383, 408
- Netflix challenge 44
- network 1 ff
- network backbone provider 15
- network models
 - acyclic networks 419, 438
 - assortative mixing 420, 421
 - Barabási–Albert model 448
 - bipartite networks 419, 433
 - cascade model 368
 - Chung–Lu model 375, 377, 423
 - clustering 368, 421, 425
 - community structure 421, 522
 - configuration model 370
 - degree-corrected stochastic block model 422, 522
 - degree correlations 420
 - directed networks 416, 433
 - duplication–divergence process 478
 - dynamic networks 424
 - Erdős–Rényi model 345
 - growing networks 434, 472
 - hierarchical networks 553, 724
 - network formation 434
 - network optimization 479
 - node copying 472, 478
 - planted partition model 543
 - Poisson random graph 343, 345, 347
 - preferential attachment 435, 448, 458
 - Price model 435
 - random graph 342, 369, 416
 - small-world model 425
 - stochastic block model 421, 522, 542
 - transitivity 368, 421, 425
 - Watts–Strogatz model 425
- network navigation 64, 718
- network optimization 479
 - airline network 479, 480, 486
 - hub-and-spoke network 479, 482, 486
 - road network 486
 - star graph 481
- network search 710
 - hierarchical model 724
 - peer-to-peer networks 42, 713
 - small-world experiment 64, 718
 - web search 35, 166, 710
- network visualization 8, 145, 219
 - and community detection 496
 - and embedding 145, 562
 - and graph Laplacian 145
 - spectral method 145
- NetworkX (software library) 220
- neural network 5, 88
 - C. elegans* 92
 - empirical measurement 90
 - functional connectivity 94
- neuron 88
- neutral fixed point 679
- news spreading 10, 607, 622
- node connectivity 138
 - algorithm for 269
- node copying 472
 - and preferential attachment 475
 - biological networks 474
 - citation networks 473
 - degree distribution 472, 476
 - duplication–divergence process 478

INDEX

- exponent 476
- mapping to preferential attachment 475
- master equation 475
- metabolic networks 477
- models **472**, 478
- power-law degree distribution **472**, 476
- protein–protein interaction networks 477
- relation to Price model 475
- node cut set 137, **139**
 - algorithm for 268, 269
 - and k -components 181
 - and robustness **181**, 262
 - minimum **139**, 181
- node disambiguation 277, **300**
- node-disjoint paths, *see* node-independent paths
- node-independent paths **137**, 269
 - algorithm for 269
 - and k -components 180
- node percolation, *see* site percolation
- nodes 1, 105
 - average degree **127**, 130
 - centrality 9, **159**
 - citation networks 4, **37**
 - copying 472
 - cut set 137, **139**
 - degree 9, **126**, 159
 - disambiguation 277, **300**
 - errors on 276
 - examples of **105**, 115
 - extraneous 276
 - failure of 569
 - food webs 5, **96**
 - groups of 177, 495
 - high degree 9, 168, **315**
 - highest degree 221, 272, 315, 328, 586, 588, 592, 601, 604, 659, 699
 - importance 9, **159**
 - Internet **15**, 19, 22
 - metabolic networks 6, **71**
 - missing 276
 - number of 106
 - percolation 569
 - power grid 27
 - removal of 464, 466, **569**, 571, 586, 593
 - social networks 4, **47**
 - values on **109**, 225
 - weighted 109
 - World Wide Web 3, **33**, 105
- non-negative matrix factorization 46
- normalized mutual information 547
- occupation probability 572
- oil pipeline network 29
- one-mode projection **116**, 178, 419
 - adjacency matrix of 118
 - cliques in **117**, 178
 - film actor network 117
 - metabolic network 72
 - rail network 29
 - weighted 117
- onion structure 556
- online dating network 49, 60
- online network 5, 49, **60**, 63
 - blogs 60, 540
 - errors in 279
 - Facebook 5, 49, **60**, 63, 363
 - journals 60
 - LinkedIn 60
 - Twitter 4, **60**
 - Usenet 49
 - weblogs 60, 540
- O notation 222
- opinion formation models 622
- optimization 479
 - airline network **479**, 480, 486
 - hub-and-spoke network **479**, 482, 486
 - modularity 498
 - road network 486
 - star graph 481
- ordered network **558**, 564
 - acyclic network 564
 - animal networks 565
 - citation networks 111, **564**
 - dominance hierarchy 58, **565**
 - faculty hiring network 565
 - food webs 564
 - friendship network 564
 - minimum violations ranking 564
 - PageRank 564
 - sports competition 565
 - statistical inference 564
- order parameter **582**, 583
- oscillators 683, 685, **701**, 709
- oscillator network 701
- out-components **136**, 308
 - giant 308
 - World Wide Web 36, 309
- out-degree 9, **130**
 - average 130
 - correlation with in-degree 274, **316**, 433
 - distribution 316, 322
 - social networks 52
 - World Wide Web 9, **316**, 317
- out-degree distribution 316
 - citation network 322
 - World Wide Web **316**, 317, 322
- overlapping communities 551
 - CFinder algorithm 551
 - edge clustering 552
 - statistical inference **551**, 553
- P2P network, *see* peer-to-peer network
- package delivery network 29
 - hub-and-spoke structure 479
- packet-switched network **15**, 25, 27
- PageRank 165
 - and rank structure 564
 - and web search **166**, 712
 - disadvantages 713
 - on a tree 213
 - parameter values 166
- pair approximation 650
- Pajek (software package) 220
- Pareto distribution 317
- partitioning 143
- patent citation network 39
- paths 131
 - algorithms for **241**, 249, 257, 262, **268**, 269
 - augmenting 264
 - directed networks 312
 - disjoint **137**, 180, 262, 268, 269

- edge-disjoint 137
- edge-independent 137, 181, 262, 268
- geodesic 132
- independent 137, 180, 262, 268, 269
- node-disjoint 137
- node-independent 137, 180, 269
- number of given length 131
- self-avoiding 131
- shortest 132, 241, 247, 249, 257, 310
- vertex-independent 137, 180, 269
- Pearson correlation
 - assortative mixing 209, 336
 - degree correlation 211, 238, 336
 - rows of adjacency matrix 197
 - similarity measure 197
- peer-to-peer network 42, 713
 - bandwidth limitations 715
 - breadth-first search on 714
 - client nodes 717
 - degree distribution 43
 - Gnutella 43
 - Napster 42
 - problems with 715
 - scale-free 43
 - searching 42, 713
 - small-world effect 715
 - supernodes 43, 716
- percolating cluster 573
- percolation 569
 - algorithms for 594
 - and epidemics 625, 627, 629, 633, 637
 - and robustness 569, 578, 580, 586, 599
 - and vaccination 570, 586, 672
 - Bethe lattice 575
 - bond percolation 571, 627
 - bootstrap percolation 179, 639
 - by degree 572, 586, 588, 601
 - clusters 573, 594, 599, 627
 - coauthorship network 600, 602
 - computational complexity 594, 598
 - computer simulation 594
 - configuration model 574, 586, 602, 603, 604, 629, 633, 637
 - disease spreading 625, 627, 629, 633, 637
 - edge percolation 571, 627
 - epidemic threshold 628
 - exponential degree
 - distribution 580, 588, 603, 605, 671
 - finite-size effects 584
 - geometric degree distribution 580, 588, 603, 605, 671
 - giant cluster 572, 574, 580, 587, 590, 599, 627, 629
 - graphical solution 575
 - immunization 570, 572, 579, 586, 592, 605, 672
 - Internet 569, 578, 600, 602
 - joint site/bond percolation 672
 - lattice 603
 - non-uniform 572, 586, 599, 601, 604
 - occupation probability 572
 - percolating cluster 573
 - phase transition 582
 - power grid 599, 601
 - power-law degree distribution 578, 579, 584, 590, 600, 602, 604, 633
 - random graph 578, 602, 604, 631
 - real-world networks 593, 599
 - regular network 602, 671
 - relabeling algorithm 596
 - road network 599, 601
 - scale-free networks 578, 579, 584, 590, 600, 602, 604, 633
 - small clusters 572, 602, 627
 - social network 600, 602
 - spanning cluster 573
 - square lattice 603
 - threshold 573, 577, 581, 590, 599, 602, 603, 604, 627, 630
 - tree 575
 - uniform 571, 574, 599
 - vaccination 570, 572, 579, 586, 592, 605, 672
- percolation threshold 573, 577, 582, 627
 - bond percolation 630
 - configuration model 577, 590, 602, 603, 604, 630
 - epidemic threshold 628
 - exponential degree
 - distribution 581, 603
 - Poisson degree distribution 578, 631
 - power-law degree distribution 578, 579, 600, 604, 633
 - random graph 578, 631
 - regular network 602, 671
 - scale-free network 578, 579, 600, 604, 633
 - sharpness 584
 - site percolation 573, 577, 581, 590, 599, 602, 603, 604
- periphery 180, 210, 555
- Perron–Frobenius theorem 160, 166, 169, 699
 - proof 161
- personal network 55
- phase transition 349, 582
 - complex contagion 644
 - configuration model 385, 390, 574, 577, 630
 - continuous 582
 - discontinuous 582, 644
 - epidemic transition 616, 628
 - finite-size effects 584
 - first-order 582, 644
 - order parameter 582, 583
 - percolation 573, 582, 584
 - random graph 349, 351, 359
 - second-order 583
 - synchronization 706
- pipeline network 29
- planar network 123
 - and coloring 125
 - approximate planarity 126, 486
 - average degree 156
 - Euler characteristic 156

INDEX

- examples 123
- faces 156
- four-color theorem 124
- Kuratowski's theorem **126**, 155
- network of regions on a map 124
- river networks 123
- road networks 124
- test for 125
- trees 123
- planted partition model 543
- plant network
 - root network 31
 - plant–pollinator network 100
- Poisson degree distribution **346**, 365, 407
 - epidemic threshold 631
 - complex contagion 642
 - configuration model **372**, 407
 - generating functions 407
 - percolation **578**, 604, 631
 - random graph **346**, 365
 - robustness 578
 - stochastic block model 422
- Poisson distribution 347
 - generating function **403**, 407
- Poisson random graph, *see* random graph
- political blog network 540
- post office delivery network 29
- power failures 28, 599, 601, 608
- power grid 27
 - degree distribution 599
- edges 27
- failure of 28, 599, 601, 608
- geographic structure 28
- nodes 27
- percolation on 599, 601
- power-law degree distribution 20, **317**, 395, 403, 435
 - and epidemics 633
 - Barabási–Albert model **450**, 455, 471
 - citation network 37, 39, 317, 322, **435**, 443
 - configuration model 372, **395**, 578, 585, 633
 - cumulative distribution **321**, 447
 - detection 319
 - exponent **318**, 329, 442, 450, 476, 578
 - generating functions 397, **403**, 408
 - Internet 20, 280, **317**, 321, 322, 578
 - metabolic network 477
 - node copying model **472**, 476
 - peer-to-peer networks 43
 - percolation **578**, 579, 584, 590, 600, 602, 604, 633
 - percolation threshold **578**, 579, 600, 604, 633
 - preferential attachment **435**, 442, 450, 461
 - Price model 442
 - robustness 578, 585, **591**, 600, 602
 - tail **326**, 404, 442
 - World Wide Web **317**, 322
- power-law distribution 317
 - betweenness centrality 331
 - centrality measures 330
 - cumulative distribution **321**, 447
 - degrees 20, 43, **317**, 372, 395, 404, 435, 442, 476, 578, 633
 - detection 319
 - diverging moments 240, **328**, 382, 397, 578, 633, 699
 - eigenvector centrality 330
 - exponent **318**, 324, 329, 442, 450, 476, 578
 - first moment 327
 - generating functions 397, **403**, 408
 - Hill estimator 324
 - Lorenz curve **328**, 338
 - maximum likelihood estimator 324
 - mean 327
 - mean square 240, 327, **328**, 382, 397, 578, 633, 699
 - mechanisms for **435**, 448, 472
 - moments 327
 - non-network examples 325
 - normalization 326
 - properties of 325
 - pure form **326**, 395, 404, 604
 - rank–frequency plot 323
 - second moment 240, 327, **328**, 382, 397, 578, 633, 699
 - tail **326**, 442
 - tests for 319
 - visualization 319
 - Yule distribution 442
- power outage 28, 608
- predator–prey interactions **96**, 108, 110, 685
- preferential attachment **435**, 448
 - addition of extra edges 459
 - and node copying 475
 - average degree **436**, 456
 - Barabási–Albert model **448**, 459
 - citation networks **435**, 443, 472
 - computer simulation **443**, 450
 - cumulative advantage 435
 - cumulative degree distribution 447
 - degree as a function of time **456**, 457
 - degree distribution **438**, 446, 450, 461, 465, 468
 - drawbacks 436, 443, 476
 - empirical evidence for 466
 - exponent **442**, 450, 461
 - generalizations 458
 - master equation **439**, 449, 451, 460, 462, 467, 488
 - models **435**, 448, 458
 - multiedges **437**, 439
 - non-linear 466
 - power-law degree distribution **442**, 450, 461
 - Price model **435**, 451, 475
 - relation to node copying process 475
 - removal of edges 461
 - removal of nodes 464, **466**
 - sublinear 466, **468**

- superlinear 472
- Yule process 435
- preflow-push algorithm 263
- Price model 435
 - acyclic networks 438
 - average degree **436**, 456
 - citation networks **435**, 443
 - computer simulation 443
 - cumulative degree distribution 447
 - degree distribution **438**, 446, 455
 - degree as a function of time **456**, 457
 - degree-zero nodes 440
 - drawbacks 436, 443, 476
 - dynamics 451
 - exponent 442
 - extensions 459, 488, 489
 - first mover effect **451**, 456
 - in-degree distribution 438
 - initial conditions 437
 - master equation **439**, 451
 - multiedges **437**, 439
 - power-law degree distribution 442
 - relation to Barabási–Albert model 449
 - relation to node copying model 475
 - time evolution **451**, 456
 - without preferential attachment 488
 - World Wide Web 438, 459
- prior probability **283**, 291, 521
- probability generating function, *see* generating function
- profile likelihood **528**, 567
- projection of bipartite network, *see* one-mode projection
- protein **76**, **81**
 - bait protein 79
 - complex 77
 - interaction network **6**, **76**
 - prey protein 79
 - structure network 87
- protein–protein interaction network **6**, **76**
 - affinity purification 80
 - bipartite representation 77
 - co-immunoprecipitation 77
 - databases 80
 - errors in **79**, 279
 - measurement 77
 - models of 477
 - mutations 478
 - node copying 477
 - S. cerevisiae* 80
 - tandem affinity purification 80
 - yeast two-hybrid screen 78
- protein structure network 87
- proximity network 59, 293
- pure power law **326**, 395, 404, 604
- questionnaires 51
- queues 245
- rail network 28
- Ramón y Cajal, Santiago 92
- Rand index 546
- random-field Ising model 142
- random graph **342**, 345, 369
 - acyclic networks 419
 - assortative mixing **364**, 420, 421
 - average component size **355**, 391
 - average degree in small components **358**, 392
 - average degree 345
 - average number of edges 344, **345**, 370
 - Bernoulli random graph **343**, 345
 - bicomponents 367
 - bipartite **419**, 433
 - bond percolation on 604, **629**
 - clustering coefficient 332, 334, **347**, 364, 366, 381
 - coinfection on 638
 - community structure 365, 421
 - complex contagion on 642
 - components **347**, 354, 355, **391**, 411, 414, 419
- configuration model 369
- cross-immunity on 637
- degree correlations 364, 420
- degree distribution **346**, 365, 370
- diameter **360**, 399
- directed 367, **416**, 433
- divergence of component sizes **359**, 395
- drawbacks 364
- dynamic networks 424
- ensemble **343**, 344
- epidemic threshold 631
- Erdős–Rényi model 345
- extensive components **348**, 354
- fit to data 520
- fixed degree distribution 371
- fixed degree sequence 370
- fixed edge probability **344**, 375
- fixed expected degrees 375
- fixed number of edges **343**, 370
- generating functions 407
- giant bicomponent 367
- giant component **347**, 354, **384**, 616
- $G(n, m)$ 343
- $G(n, p)$ 344
- graphical solution **351**, 389
- hierarchical 553
- homophily model 421
- k -components 367
- large size limit 344, 346, 349, 354, 356, 362, 364
- largest component 348
- largest eigenvalue 708
- locally tree-like 122, 353, 356, **382**, 383, 386, 391, 654
- loops in **366**, 421
- multiedges **371**, 373
- number of edges 343, **345**, 370
- path lengths 360, 399
- percolation on **578**, 602, 604, 631
- percolation threshold **578**, 631
- phase transition **349**, 351, 359, 385

INDEX

- Poisson degree distribution **346, 365**
- problems with 364
- regular network 429
- robustness 578
- self-edges **371, 374**
- simple network 343
- SIR model on 631
- small components **355, 391, 411, 414**
- small-world effect **360, 363, 401**
- temporal networks 424
- transitivity **347, 364, 366, 368, 421**
- tree-like components **356, 382, 391**
- triangles in **347, 364, 366, 368, 421**
- two-mode networks **419, 433**
- with arbitrary degree distribution 369
- with assortativity 420, 421
- with community structure 421
- with degree correlations 420
- with given expected degrees 375
- with transitivity 368, 421
- random walk 147
 - and graph Laplacian 147
 - and node degree 68, **148**
 - betweenness 177
 - community detection **515, 553**
 - entropy 518
 - sampling 67
 - search strategy 730
 - self-avoiding 147
- random-walk betweenness 177
- random-walk sampling **67, 148**
- random-walk search 730
- rank–frequency plot 323
- rank structure **558, 564**
 - acyclic networks 564
 - animal networks 565
 - citation networks 564
 - dominance hierarchies 58, 565
 - faculty hiring network 565
 - food webs 564
 - friendship networks 564
 - minimum violations ranking 564
 - PageRank 564
 - sports competition 565
 - statistical inference 564
- reality mining study 293
- reciprocated edges 189
- reciprocity 189
 - calculation of 274
 - email network 190
 - friendship network 189
 - World Wide Web 190
- recommender network **44, 115**
- recommender system 44
- recovered state 612
- recovery rate 613
- red-black tree 122, 234
- redundancy 186
 - and local clustering coefficient 188
- Reed–Frost model 672
- regional ISP 16
- regular equivalence **194, 198**
 - measures of 198
- regular graph, *see* regular network
- regular network 128
 - centrality 211
 - configuration model 429, 602, 671, 673
 - dynamical systems on 706
 - eigenvector centrality 211
 - epidemics on 671, 673
 - Katz centrality 211
 - largest eigenvalue 708
 - lattice 128
 - percolation on 602, 671
 - random 429
 - SI model on 673
 - SIR model on 671
- regulatory network, *see* genetic regulatory network
- reinfection 617
 - SIRS model 619
 - SIS model **617, 667**
- removed state 613
- Rényi, Alfréd 345, 350
- repelling fixed point **679, 681, 688**
- resilience, *see* robustness
- resistor network 148
- respondent-driven sampling **68, 593**
- reverse small-world experiment **65, 723**
- rich-get-richer effect 435
- Riemann zeta function 326, 396, 404
 - incomplete 326
- river network 29
 - geographic structure 30
 - planarity 123
 - Russia 28
 - tree **31, 121, 123**
- RNA structure network 87
- road network **28, 561**
 - embedding 561
 - latent-space structure 561
 - model for 486
 - optimization 486
 - percolation on **599, 601**
 - planarity 124
- robustness **181, 569, 578, 586**
 - and connectivity **181, 262**
 - and cut sets **181, 262**
 - and epidemics **579, 590**
 - and independent paths **181, 262**
 - and k -components 181
 - configuration model **578, 583, 585, 590**
 - exponential degree distribution 580, 588
 - independent paths **181, 262**
 - Internet 181, **569, 571, 573, 578, 592, 600, 602**
 - node removal **569, 578, 580, 586, 599**
 - percolation **569, 578, 580, 586, 599**
 - Poisson degree distribution 578
 - power-law degree distribution 578, 585, **591, 600, 602**
 - random graph 578
 - real-world networks 599

- scale-free networks 578, 585, 590, 600, 602
- social networks 570, 600
- targeted node removal 586, 591, 601
- root network 31
- root node 121
- rooted tree 121
- router 15
 - failure 181, 569, 571
 - routing tables 21
- Routeviews Project 22
- routing table 21
- rumors 10, 607, 622, 693
- running time 221
 - adjacency list operations 232, 234
 - adjacency matrix operations 226
 - augmenting path algorithm 266
 - betweenness centrality 255, 257
 - breadth-first search 222, 246
 - closeness centrality 249
 - clustering coefficient 239
 - diameter 249
 - Dijkstra's algorithm 260
 - hierarchical clustering 536
 - Louvain algorithm 512
 - InfoMap algorithm 520
 - modularity maximization 504, 507, 512
 - on sparse networks 223, 227, 232
 - percolation algorithms 594, 598
 - reciprocity 274
 - shortest distance 246, 249, 260
 - shortest path 246, 249, 251, 260
- Saccharomyces cerevisiae* 80
- saddle point 681, 688
- Salton cosine similarity 196
- scale-free network 317, 435
 - average distance on 312
 - Barabási–Albert model 448, 450
 - biological networks 477
 - citation networks 37, 39, 317, 322, 435, 443
 - configuration model 372, 395, 578, 585, 633
 - core 312
 - cumulative degree distribution 321, 447
 - degree distribution 20, 43, 240, 317, 441, 450
 - diameter 312
 - duplication–divergence model 472
 - epidemics on 633
 - epidemic threshold 633
 - exponent 318, 321, 324, 329, 395, 442, 450, 461, 465, 476, 578
 - fragility 590, 602
 - generating functions 397, 403, 408
 - giant cluster 584, 590, 600, 602, 633
 - giant component 396
 - immunization 579, 592
 - Internet 20, 280, 317, 321, 322, 578
 - Lorenz curve 328
 - mean square degree 240, 328, 382, 397, 578, 633, 699
 - metabolic networks 477
 - multiedges 374
 - node copying model 472
 - non-uniform percolation on 590, 602
 - peer-to-peer networks 43
 - percolation on 578, 579, 584, 590, 600, 602, 604, 633
 - percolation threshold 578, 579, 600, 604, 633
 - preferential attachment models 435, 442, 448, 458
 - Price model 435
 - real-world networks 317
 - robustness 578, 585, 590, 600, 602
 - shortest paths 312
 - SIR model on 633
 - vaccination 579, 592
 - World Wide Web 317, 322
- S. cerevisiae* 80
- Science Citation Index 37
- scientometrics 37
- Scopus 37
- search 710
 - breadth-first 34, 241
 - engine 35, 166, 710
 - greedy algorithm 719
 - hierarchical model 724
 - Kleinberg model 718
 - peer-to-peer networks 42, 713
 - random-walk search 730
 - small-world experiment 64, 718
 - web search 35, 166, 710
- second neighbors 383, 408
- second-order phase transition 583
- seed dispersal network 100
- SEIR model 620
- self-avoiding path 131, 133
- self-avoiding random walk 147
- self-avoiding walk 131
- self-edge 106
 - acyclic networks 111
 - adjacency list representation 230
 - adjacency matrix representation 107, 110
 - configuration model 371, 374
 - directed networks 110, 130
 - multiple 108
 - random graphs 371, 374
- self-loop, *see* self-edge
- sewerage network 29
- sexual contact network 49, 53, 67
 - assortative mixing 216
 - bipartite 116
 - disassortative mixing 202, 205
 - modularity 205
- Shannon source coding theorem 518
- shortcut 426
- shortest distance 10, 132
 - algorithm for 241, 247, 257
 - all pairs of nodes 249

INDEX

- average 170, **311**, 363
- computational complexity **246**, 249, 260
- diameter **133**, 249
- Dijkstra's algorithm 258
- harmonic mean 172
- infinite 133
- log n behavior 312
- random graph **360**, 399
- real-world networks **310**, 363
- scale-free network 312
- small-world effect **132**, 310, 363
- weighted networks 257
- shortest path **132**, 241, 257, 310
 - absence of loops in 133
 - algorithms for **241**, 247, 249, 257
 - computational complexity 246, 249, 251, 260
 - Dijkstra's algorithm **258**, 261
 - directed networks 312
 - longest **133**, 249
 - multiple **133**, 251
 - on a tree **251**, 253, 261
 - overlapping **133**, 251
 - real-world networks **310**, 363
 - self-avoiding 131, **133**
 - uniqueness **133**, 251
 - weighted networks 257
- shortest-path tree **251**, 261
 - betweenness centrality 253
 - Dijkstra's algorithm 261
- signed edges 84, 89, 109, **190**, 258
- signed network 190
 - clusterability **192**, 216
 - frustration 191
 - loops in 190
 - structural balance 192
 - triangles in 190
- similarity 194
 - and community detection 534
 - and disambiguation 300
 - and entity resolution 300
 - and link prediction 299
 - between groups of nodes 535
 - bibliographic coupling 39
 - cocitation 39
 - correlation coefficient 197
 - cosine similarity **196**, 299, 300, 537
 - documents 45
 - Euclidean distance 197
 - Hamming distance 197
 - hierarchical clustering 534
 - Jaccard coefficient **197**, 299
 - Katz similarity 200
 - latent semantic analysis 46
 - measures of 195, 198
 - number of common neighbors 195
 - Pearson correlation 197
 - regular equivalence **194**, 198
 - structural equivalence **194**, 195
- SI model **609**, 625
 - and eigenvector centrality 648
 - as a dynamical system **646**, 676, 678, 685, 688
 - configuration model **655**, 673
 - degree-based approximation 654
 - early-time behavior 611, **647**
 - fixed points 678, 688
 - fully mixed **611**, 626, 646, 676
 - initial conditions 647
 - late-time behavior 625, 648, 657
 - logistic growth 611
 - moment closure method 651
 - on a network **625**, 646, 650, 654
 - outbreak size 625
 - pair approximation 650
 - regular network 673
 - short-time behavior 611, **647**
 - solution **611**, 647, 650, 654
 - symmetric fixed point 689
 - time evolution 611, **646**, 650, 654
- simple graph, *see* simple network
- simple network 106
 - adjacency matrix **106**, 131
 - density 128
 - disassortative mixing 337
 - maximum number of edges 128
 - random graph 343
 - single-linkage clustering **535**, 537
 - singular value decomposition 46
 - sink food web 99
- SIR model **612**, 626
 - and bond percolation **627**, 629
 - and eigenvector centrality 661
 - basic reproduction number 617
 - configuration model **629**, 662, 671
 - degree-based approximation 662
 - early-time behavior 615, 660, 666
 - epidemic outbreaks 616, 627, 629, 631, 661, 665
 - epidemic threshold 616, 617, 630, 632, 661, 666
 - fully mixed **613**, 624, 626, 629, 632
 - initial conditions **615**, 660
 - late-time behavior 626, 665
 - outbreak size **626**, 629, 631, 665
 - power-law degree distribution 633
 - probability of epidemic 631
 - random graph 631
 - regular networks 671
 - scale-free networks 633
 - small outbreaks 672
 - solution **614**, 629, 660, 663
 - time evolution 614, **660**, 662
 - transmission probability 613, 624, **626**
 - transmission rate 613, 624
 - vaccination 672
- SIRS model 619
 - fully mixed 620
- SIS model **617**, 667
 - and eigenvector centrality 668
 - basic reproduction number 618
 - degree-based approximation 669
 - early-time behavior **667**, 669
 - endemic state 670

- epidemic threshold 618, 668, 670
- fully mixed 617
- initial conditions 618, 667, 669
- late-time behavior 618, 668, 670
- logistic growth 618
- solution 618, 667, 669
- time evolution 617, **667**, 669
- site 105
- site/bond percolation 672
- site percolation **569**, 571
 - algorithms for 594
 - and robustness **569**, 578, 580, 586, 599
 - Bethe lattice 575
 - by degree **572**, 586, 588, 601
 - clusters **573**, 594, 599
 - coauthorship network **600**, 602
 - computational complexity **594**, 598
 - computer simulation 594
 - configuration model **574**, 586, 602, 603, 604
 - exponential degree
 - distribution **580**, 588, 603, 605
 - finite-size effects 584
 - geometric degree distribution **580**, 588, 603, 605
 - giant cluster **572**, 574, 580, 587, 590, 599
 - graphical solution 575
 - immunization **570**, 572, 579, 586, 592, 605, 672
 - Internet **569**, 571, 573, 578, 592, 600, 602
 - joint site/bond percolation 672
 - non-uniform 572, **586**, 599, 601, 604
 - occupation probability 572
 - percolation threshold **573**, 577, 581, 590, 599, 602, 603, 604
 - phase transition **573**, 582, 584
 - power grid **599**, 601
 - Poisson degree distribution
 - 578
 - power-law degree distribution
 - 578**, 579, 584, 590, 600, 602, 604
 - random graphs **574**, 587
 - random removal of nodes **571**, 574, 599
 - real-world networks **593**, 599
 - regular network 602
 - relabeling algorithm 596
 - road network **599**, 601
 - scale-free networks **578**, 579, 584, 590, 600, 602, 604
 - small clusters **572**, 602
 - social network **600**, 602
 - spanning cluster 573
 - uniform 571, 574
 - vaccination **570**, 572, 579, 586, 592, 605, 672
- Six Degrees of Kevin Bacon 61
- six degrees of separation **10**, 63, 310, 718
- small components **306**, 308
 - average degree in 358, 392
 - average size **356**, 359, 391, 394
 - configuration model **391**, 411, 414
 - degree distribution in 366, 393
 - directed networks 308
 - diverging size 359, 395
 - exponential degree
 - distribution 416
 - generating functions **411**, 414
 - number of 356
 - random graph **355**, 391, 411, 414, 432
 - size distribution 356, **414**, 432
 - trees **356**, 391
- small-world effect **10**, **62**, 310
 - and disease spreading 311
 - and rumors **10**, 311
 - and search 715, **718**
 - configuration model 401
 - Facebook network **63**, 312, 363
 - Internet **10**, 311
 - message passing experiments
 - 61, **62**, 310, 313, 363, 718, 723
 - models of 312, **360**, 363, 401, 426
 - peer-to-peer networks 715
 - random graph **360**, 363, 401
 - real-world networks 63, **310**, 363
 - small-world model **426**, 718
 - small-world experiment 61, **62**, 310, 718
 - and network search 64, **718**, 723
 - email version 64
 - errors in **63**, 64
 - funneling **64**, 313
 - hierarchical model 723
 - INDEX experiments 65
 - Kleinberg model **718**, 720
 - models of **718**, 720, 723
 - problems with 63
 - response rate **63**, 64
 - reverse small-world
 - experiment **65**, 723
 - small-world model 718
 - small-world model 425
 - Kleinberg variant 718, **720**
 - shortcuts 426
 - snowball sampling 65
 - biases 66
 - social contagion 607, **622**, 639
 - social network 4, **47**
 - actors 47, 61, 106
 - Add Health study **54**, 201
 - affiliation networks 49, **60**, 114
 - alters 55
 - animals 49, **58**, 279, 565
 - animosity 109, **190**
 - ants 58
 - archival data 49, **58**, 60, 279
 - assortative mixing **201**, 209, **335**
 - baboons 58
 - bipartite 61
 - bison 58
 - board of directors 49, **61**, 114
 - brokers 176
 - businesspeople 49, **61**
 - cellphones 59, 293

INDEX

- clustering coefficient **185**, 333
- coauthorship 49, **61**, 185, 333, 380, 495
- community structure 201, 206, 495, 503, 539
- dating 49, 60
- deer 58
- degree correlations 336
- degree distribution 52, 56
- directed networks **52**, 185
- direct observation **57**, 279
- disassortative mixing **202**, 205
- disease spreading 53, 67, **607**
- doctors 53
- dolphins 58, **539**
- drug users 49
- dynamic 60, 120
- edges 4, **47**, 48, 106
- ego 55
- ego-centered 55
- email networks **59**, 64, 185, 190, 279
- errors in 52, 54, 63, 66, 67, 276, **278**, 293
- ethnographic studies **57**, 279
- Facebook 5, 49, **60**, 63, 279, 312, 363
- face-to-face interactions **57**, 59, 293
- film actors **61**, 115, 117, 171, 176, 185, 306
- fixed choice studies **52**, 279
- Florentine families 58
- free choice studies **52**, 56, 279
- friendship 5, 48, **51**
- geographic structure 24, 59, 209, 559, 724
- historical 48, 49, 58
- homophily 54, 55, **201**, 206, 335
- hubs in 52
- in-degree 52
- instant messaging 59
- interviews 51
- intermarriage network 58
- kangaroos 58
- karate club network 5, 57, **503**, 507, 537, **539**
- LinkedIn 60
- measurement 49, **51**, 55, 57, 58, 60, 62, 276, 278, 293
- missing data 52, 54, 63, **278**
- mobile phones 59, 293
- modularity **205**, 216
- monkeys 58
- multilayer 54, **120**
- multiplex 54, **120**
- musicians 49
- name generators **51**, 201
- National Longitudinal Study of Adolescent Health **54**, 201
- nodes 4, **47**, 106
- online 4, 49, **60**, 63, 279, 312, 363
- online dating 60
- online diaries 60
- out-degree 52
- percolation on **600**, 602
- personal networks 55
- proximity networks 59, 293
- questionnaires 51
- respondent-driven sampling **68**, 593
- robustness 570, 600
- rumor spreading 10, 607, **622**, 693
- schoolchildren 49, **51**, 201, 205, 206, 279, 496
- scientists 49, **61**, 185, 333, 380, 495
- sexual contacts 49, **53**, 67, 116, 202, 205, 216
- students 5, 57, 293, **503**
- sociometric studies 55
- Southern Women Study **49**, 60
- sparse 129
- sports competition networks 565
- strength of edges **54**, 108
- surveys **51**, 55, 65, 278
- telephone surveys **51**, 55, 66
- temporal 60, 120
- terrorists 49
- text messaging 59
- third-party data 49, **58**, 60, 279
- ties 47, 106
- time-resolved 60, 120
- transitivity **185**, 333
- types of edges 53, 120
- university students 5, 49, 57, 293, 503, 507, 537, 539
- Usenet 49
- weighted **54**, 108
- windsurfers 57
- wolves 58
- social network analysis 47, 158
- Social Science Citation Index 37
- sociogram 48
- sociometric superstar **64**, 313
- software call graph 25
- software 218
 - call graphs 25
 - for network analysis **219**, 236, 498, 511
 - Gephi 220
 - Graphviz 220
 - InFlow 220
 - JUNG 220
 - network visualization 219
 - NetworkX 220
 - Pajek 220
 - UCINET 176, **220**, 228
 - Visone 220
 - yEd 220
- source coding theorem 518
- source food web 99
- Southern Women Study 49
 - bipartite network 49, 60
- spanning percolation cluster 573
- sparse network 128
 - adjacency matrix of 227
 - algorithms on **223**, 232, 247, 249, 257, 261, 266, 508, 512, 520, 599
 - and computational complexity **223**, 227, 232
 - constant average degree **129**, 223
 - extremely sparse **129**, 223
 - Internet 129

- running time of algorithms
 - 223, 227, 232
- social networks 129
- World Wide Web 129
- sparsification 150
- spectral partitioning 145
- spectra of networks 698
 - adjacency matrix 160, 164, 661, 690, 697, **698**, 708
 - graph Laplacian **150**, 701
- sports competition network 565
- square lattice
 - diameter 337
 - eigenvalues and eigenvectors 708
 - percolation on 603
- star graph 481
 - airline network 481
 - betweenness centrality 215
 - largest eigenvalue 154
- static web pages 35
- statistical inference
 - community detection 520, **522**, 551, 553
 - core-periphery structure 558
 - embedding 563
 - error estimation 281, 285
 - normal distribution 281
 - hierarchical structure **553**, 564
 - latent-space structure 563
 - overlapping communities **551**, 553
 - rank structure 564
- stochastic block model **421**, 522, 542
 - assortative mixing 421
 - benchmark networks **542**, 543, 550
 - community detection 423, **522**, 542
 - core-periphery structure 558
 - degree-corrected **422**, 522, 543
 - degree distribution 422
 - disassortative mixing 422
 - homophily 421
 - likelihood **523**, 526
 - link prediction 299
 - mixed membership 552
 - profile likelihood **528**, 567
 - stochastic dynamical system 675
 - stratified network 206, 558
 - stretched exponential 465, 470
 - compressed exponential 471
 - degree distribution 465, **470**
 - strongly connected component **135**, 308
 - acyclic networks 309
 - eigenvector centrality in 162
 - Katz centrality in 163
 - largest 308
 - small 309
 - World Wide Web 135, **308**
 - structural balance **190**, 192
 - Davis's balance theorem **193**, 216
 - Harary balance theorem **192**, 216
 - structural equivalence 194
 - measures of 195
 - structural holes 186
 - subnet 20
 - supernode 43, **716**
 - superpeer 716
 - surveys 51
 - Add Health study **54**, 201
 - design of 51
 - ego-centered networks 55
 - errors in 54, 66, **278**
 - fixed choice **52**, 279
 - free choice **52**, 56, 279
 - General Social Survey 55
 - name generators **51**, 201
 - respondent-driven sampling **68**, 593
 - snowball sampling 65
 - susceptible-infected model, *see* SI model
 - susceptible-infected-removed model, *see* SIR model
 - susceptible-infected-susceptible model, *see* SIS model
 - susceptible state 609
 - symmetric fixed point **688**, 693
 - synapse 89
 - synchronization 701
 - Floquet theory 706
 - master stability function 704
 - phase transition 706
 - stability 703
- synthetic network 542
- tandem affinity purification 80
 - errors in 280
- technological network 1, **14**
 - airline network **28**, 479, 480
 - delivery networks 29
 - distribution networks 29
 - electrical networks 25, 27, 148, 599
 - Internet 1, **15**
 - measurement 17
 - pipelines 29
 - power grid **27**, 599
 - rail networks 28
 - resistor networks 148
 - road networks **28**, 124, 486, 561
 - software call graphs 25
 - telephone network 25
 - transportation networks **28**, 118, 121, 479, 480
- telephone call graph 59
- telephone network 25
 - call graph 59
 - circuit-switched 25
 - geographic structure 27, 59
 - mobile phones 59, 293
- telephone surveys **51**, 55, 66
- temporal network 118, **120**
 - multilayer network 120
 - Price model 451
 - random graph model 424
 - social networks 60, 120
- text messaging network 59
- ties 47, 106
- time complexity, *see* running time
- topology 7
- traceroute **17**, 280
- train network 28
- transitive triples 183, **185**
- transitivity **183**, 332
 - cliques 183
 - coauthorship networks 185, 333

INDEX

- configuration model 332, **381**
- directed networks 185
- email network 185
- film actor network 185
- food webs 334
- Internet 333
- local clustering coefficient **186**, 334
- models of 368, 421, 425
- random graph 334, **347**, 364, 366, 368, 421
- real-world networks 332
- small-world model 425
- social networks **185**, 333
- World Wide Web 334
- transmission probability 626
 - SIR model 613, 624, **626**
- transmission rate 613, 624
- transportation network 28
 - airline network **28**, 479, 480
 - delivery networks 29
 - geographic structure 28
 - interstate highway network 561
 - models of 479
 - multilayer **118**, 121
 - optimization 479
 - rail networks 28
 - river networks 28
 - road networks **28**, 124, 486, 561
 - weighted 109
- trapezium rule 469
- trapezoidal rule **469**, 721
- tree 121
 - absence of loops in 31, **121**, 184, 251, 356
 - average degree 154
 - AVL tree 122, 234
 - Bethe lattice 122, 575
 - betweenness centrality 213
 - binary 553, 725
 - Cayley tree 122, 338, 575
 - closeness centrality 213
 - clustering coefficient 184
 - configuration model **382**, 383, 386, 391, 654
 - data structures 122, 234
 - dendrogram 122, **531**, 553, 724
 - directed 121
 - hierarchical random graph 553
 - hierarchical structure 122, 537, 553, 724
 - leaf nodes **121**, 255, 531, 553
 - locally tree-like network 353, **382**, 383, 386
 - minimum spanning tree 122
 - number of edges **123**, 481
 - optimized networks 481
 - PageRank 213
 - percolation on 575
 - planarity 123
 - random graphs 122, 353, **356**, 382, 383, 386, 391, 654
 - red-black tree 122, 234
 - river networks **31**, 121, 123
 - rooted 121
 - root node 121
 - shortest-path tree **251**, 253, 261
 - small components **356**, 391
 - star graph 481
- triadic closure 333
- triangles **183**, 239, 332
 - closed triads **183**, 333
 - clustering coefficient **183**, 239, 332
 - community detection 532
 - in random graphs 366, 368, 421
 - in signed networks 190
 - models of 368, **421**
 - transitive triples 183, **185**
- tricomponent 181
- tripartite network 72
- trophic level **98**, 156, 564
- trophic species 96
- two-hybrid screen 78
 - advantages and disadvantages 79
 - errors in **79**, 280
- two-mode network 61, **115**
 - affiliation network **61**, 114
 - average degree 155
 - coauthorship network 61
 - film actor network **61**, 115, 117
 - incidence matrix 116
 - models for **419**, 433
 - one-mode projection 29, 72, **116**, 178, 419
 - random graph model **419**, 433
 - sexual contact 116
 - social network 61
 - Southern Women Study **49**, 60
- UCINET (software package) 176, **220**, 228
- undirected network 105
 - adjacency list 229
 - adjacency matrix **106**, 226
 - average degree 127
 - betweenness centrality 173
 - clustering coefficient 183
 - components **133**, 180, 306
 - connectivity 138
 - degree 9, **126**
 - degree distribution 313
 - degree sequence 314
 - eigenvector centrality 159
 - excess degree 380
 - excess degree distribution 377, **381**, **386**, 407, 575
 - friendship networks 52
 - k*-components 180
 - loops in 132
 - paths on **131**, 132, 134, 137, 180, 249, 263, 310
 - self-edges **107**, 110
 - sexual contact networks 53
 - social networks 52, 53
 - transitivity 183
- URL 34
- Usenet 49
- vaccination 570
 - acquaintance immunization **592**, 605
 - and network robustness 579
 - by degree 572, **586**, 590, 592
 - herd immunity **571**, 579
 - network robustness 579
 - non-uniform 572, **586**, 590, 592, 605
 - percolation theory **570**, 572, 579, 586, 592, 605, 672

- scale-free networks **579**, 591
 - targeted **572**, 586, 590, 592, 605
- valued network, *see* weighted network
- vertex 1, 105
- vertex-disjoint paths, *see* vertex-independent paths
- vertex-independent paths **137**, 269
 - algorithm for 269
 - and k -components 180
- vertices 1, 105
- Visone (software package) 220
- visualization 8, 145, **219**
 - acyclic networks 112
 - algebraic connectivity 147
 - algorithms 145
 - and community detection 496
 - and embedding **145**, 562
 - and graph Laplacian 145
 - of power laws 319
 - software 219
 - spectral method 145
- voter model 622
- walks 131
 - number of a given length 131
 - random **147**, 177, 515, 553, 730
 - self-avoiding **131**, 133, 137
- Watts–Strogatz clustering
 - coefficient **188**, 334
- Watts–Strogatz model 425
 - Kleinberg variant 718, **720**
- weakly connected component **135**, 308
 - giant **308**, 367, 433
- web crawler **33**, 307, 710
 - biases 35
 - breadth-first search 34
 - citations 37
 - operation of 33
 - software 36
- weblogs 60
 - political 540
- web pages 3, **32**
 - dynamically generated **35**, 280, 710
 - hits on 325
 - in-degree 9
 - number of **33**, 35, 280
 - out-degree 9
 - PageRank **166**, 712
 - ranking 162, **166**, 170, 712
 - reachability **35**, 135, 281, 308
 - static 35
 - URLs 34
- web search 35, 166, **710**
 - breadth-first search 34
 - crawlers **33**, 307, 710
 - Google 35, 166, 710
 - PageRank **166**, 712
 - problems with 35, 713
- web, *see* World Wide Web
- weighted edges 108
 - adjacency matrix
 - representation **108**, 226
 - errors on 277
 - examples 108
 - negative weights **109**, 258
- weighted network 108
 - adjacency matrix
 - representation **108**, 226
 - bibliographic coupling
 - network 39
 - bipartite 116
 - cocitation networks 39
 - computer representation 235
 - cut sets 141
 - Dijkstra’s algorithm 258
 - equivalence to multigraphs **109**, 141
 - errors in 277
 - food webs **100**, 108
 - graph Laplacian 143
 - Internet 108
 - max-flow/min-cut theorem 141
 - maximum flow 141
 - minimum cut 141
 - negative weights **109**, 258
 - one-mode projections 117
 - shortest distances on 257
 - shortest paths 257
 - social networks **54**, 108
 - transportation networks 109
 - weights on nodes 109
- Westlaw 41
- World Wide Web 3, **32**
 - addition of edges **458**, 459
 - assortative mixing 203
 - blogs **60**, 540
 - bow tie diagram 309
 - clustering coefficient 334
 - co-links 189
 - community structure 10, 177, 496, 540
 - components 36, 134, **308**
 - cumulative degree distribution 322
 - degree distribution 317, **322**, 329
 - directed network 4, **33**, 35, 110
 - disappearance of edges 458
 - disappearance of nodes 458
 - dynamic pages **35**, 280, 710
 - edges 3, **33**
 - errors in data 35, **280**
 - exponent 329
 - giant components 36, **308**
 - giant out-component 36, 309
 - hyperlinks 3, **33**, 105, 135, 280
 - in-components 281, **309**
 - in-degree 9, 316, 317, **322**, 329
 - in-degree distribution 316, 317, **322**, 329
 - in/out-degree correlation 433
 - largest component 36, **308**
 - measurement of 33
 - models of **435**, 438, 458, 459, 466
 - nodes 3, **33**
 - number of pages **33**, 35, 280
 - out-components 36, 309
 - out-degree **316**, 317
 - out-degree distribution **316**, 317, 322
 - pages 3, **33**, 35
 - power-law degree distribution **317**, 322
 - Price model **435**, 438
 - reachable pages **35**, 135, 281, 308
 - reciprocity 190

INDEX

- scale-free network **317**, 322
- searching 166, **710**
- size **33**, 35, 280
- sparse network 129
- strongly connected
 - components 135, **308**
- transitivity 334
- unreachable pages **35**, 135, 281, 308
- weakly connected components
 - 135
- weblogs **60**, 540
- web search 166, **710**
- yeast protein–protein interaction
 - network 80
- yeast two-hybrid screen 78
 - advantages and disadvantages 79
 - errors in **79**, 280
- yEd (software package) 220
- Yule distribution 442
- Yule process 435
- Zachary, Wayne **5**, 539
- zeta function 326, 396, 404
 - incomplete 326
- Zipf’s law 317