



Die Schwierigkeiten bei der Modellierung von Epidemien

Petra Berenbrink^(✉)

Fachbereich Informatik, MIN-Fakultät, Universität Hamburg, Hamburg,
Deutschland

berenbrink@informatik.uni-hamburg.de

Schlüsselwörter: Modellierung · Graphen · Soziale Netzwerke

1 Grundlagen

Damit sich Infektionen ausbreiten können, müssen Menschen miteinander in Kontakt stehen, also eine Art Netzwerk bilden. Bei Modellierungen solcher sozialen Netzwerke und auch in vielen anderen Anwendungen spielen sogenannte Graphen eine große Rolle. Ein Graph ist ein sehr bekanntes mathematisches Modell und besteht aus Knoten und Kanten (Abb. 1). Die Kanten verbinden jeweils zwei der Knoten. Falls der Graph ein soziales Netzwerk modellieren soll, stehen die Knoten des Graphs für die einzelnen Personen. Personen, die miteinander bekannt sind, werden mit einer Kante verbunden.

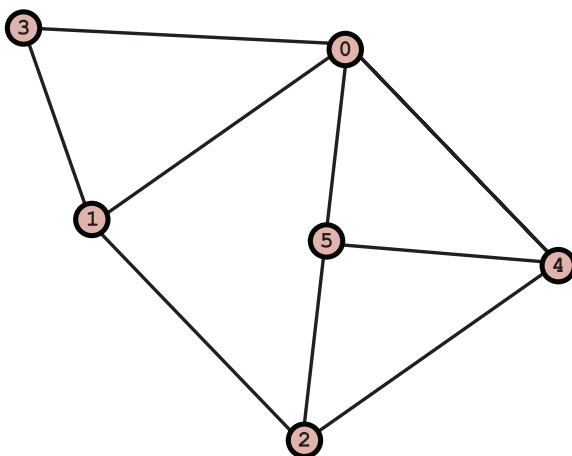


Abb. 1. Ein einfacher Graph

Das wohl bekannteste Modell für die Übertragung von Krankheiten ist das SIR-Modell, bei dem die Personen in drei Gruppen eingeteilt werden: susceptible (anfällig), infectious (infektiös) and recovered (wieder gesund). Die Personen werden

wieder durch Knoten modelliert, in der einfachsten Variante des Modells sind alle Personen miteinander verbunden. Den Prozess kann man am besten so erklären: In jedem Zeitschritt wird einer der Knoten ausgewählt. Wenn der Knoten für eine infizierte Person steht, wird ein zweiter Knoten ausgewählt. Falls der zweite Knoten eine anfällige Person repräsentiert, geht er mit einer gewissen Wahrscheinlichkeit in den infizierten Zustand über. Falls der initial ausgewählte Knoten infiziert ist, geht er mit einer gewissen Wahrscheinlichkeit in den wieder gesunden Zustand über.

In der Literatur sind sehr viele Varianten dieses Prozesses untersucht worden, die sich darin unterscheiden, wie die Knoten ausgewählt werden. Wichtig ist dabei, dass die Wahrscheinlichkeit einer Neuinfektion von der Anzahl der infizierten und der anfälligen Personen abhängt. Bei dem eben definierten und ähnlichen Modellen geht man davon aus, dass jede Person jede andere treffen und somit infizieren kann. Außerdem sind alle Personen genau gleich in dem Sinne, dass es keine Personen gibt, die besonders viele Kontakte haben oder auch besonders wenige. Mathematiker bezeichnen den zugrunde liegenden Graphen, bei dem jeder Knoten mit jedem anderen Knoten verbunden ist, als kompletten Graphen oder auch als Clique. Das gerade beschriebene Modell geht also davon aus, dass alle Personen sich genau gleich verhalten, und dass jeder Mensch jede andere Person mit der gleichen Wahrscheinlichkeit trifft.

In realistischeren Modellen sind die Personen durch einen nicht vollständigen Graphen miteinander verbunden. Jede Person hat nur eine Kante zu den Personen, zu denen Kontakte bestehen, etwa zu Arbeitskollegen, Freunden oder Familienmitgliedern. Die resultierenden sozialen Netzwerke zeichnen sich dadurch aus, dass die Knotengrade (Anzahl der Kontakte eines Knotens, auch Nachbarn genannt) eine sogenannte Powerlaw-Verteilung haben [1]. Somit haben einige wenige Knoten einen enorm großen Grad, während die meisten anderen Knoten einen recht kleinen Grad besitzen. Außerdem enthalten solche Netzwerke auch viele Cluster: kleine Knotenmengen, die sehr viele Kanten enthalten. So ein Cluster kann beispielsweise aus Arbeitskollegen, Freundeskreisen oder Mitgliedern eines Sportvereins bestehen. Es gibt auch sehr viele Dreiecke, also drei Knoten, die jeweils miteinander verbunden sind. Diese Dreiecke modellieren gemeinsame Freunde der Personen, die durch die Knoten modelliert sind. Solche Powerlaw-Netzwerke treten in sehr vielen Gebieten auf, beispielsweise als Modell für Internetverbindungen [2] oder auch für Kooperationen im wissenschaftlichen Bereich [3].

Wichtig ist nun, dass Ausbreitungsprozesse auf solchen Graphen oft sehr ähnlich aussehen. Letztendlich ist es egal, ob Krankheiten, Meinungen, oder Nachrichten und Gerüchte verbreitet werden. Die Knoten können entweder infiziert sein, wenn es um die Verbreitung von Krankheiten geht, oder informiert – sie kennen also die Nachricht oder das Gerücht. Infizierte oder informierte Knoten interagieren mit ihren Nachbarn: Das sind die Knoten, zu denen sie mit einer Kante verbunden sind und die dann später ebenfalls infiziert beziehungsweise informiert sind. Die Analysemethoden für solche zufälligen Prozesse sind sich somit sehr ähnlich. Im Bereich der Informatik wird schon lange analysiert wie sich Nachrichten verbreiten, ein Prozess, der in dem Zusammenhang auch Broadcasting genannt wird.

2 Theoretische Analyse epidemiologischer Prozesse

Das Ziel der mathematischen Analyse von epidemiologischen Prozessen ist es, diese besser zu verstehen. Mathematisch können epidemiologische Prozesse als sogenannte Markov-Kette definiert werden. Diese kann eine meistens endliche Menge von Zuständen annehmen, die man sich am besten wieder als Graph vorstellt. Die Knoten sind die Zustände. In unserem Fall sieht ein Zustand so aus: jeder Knoten ist entweder infiziert, infektiös oder wieder gesund. Zwei Knoten sind mit einer Kante verbunden, wenn von einem Zustand direkt in einen anderen übergegangen werden kann. Das sind Zustände, die sich nur in einem Knoten unterscheiden. Eine sogenannte Übergangsfunktion beschreibt, mit welcher Wahrscheinlichkeit der Prozess von einem Knoten auf einen der Nachbarknoten übergeht.

Oftmals werden Markov-Ketten mit der mean-field approximation analysiert [4]. Dabei wird von der Annahme ausgegangen, dass der Zustand eines Knotens unabhängig ist vom Zustand seiner Nachbarn. Vereinfacht ausgedrückt heißt das: Alle Knoten haben den gleichen Anteil an infizierten und nicht infizierten Nachbarn, welcher sich aus dem Quotienten der insgesamt infizierten und nicht infizierten Knoten ergibt. Eine weitere Analysemethode ist die gradbasierte Meanfield-Analyse. Hier geht man davon aus, dass sich alle Knoten mit dem gleichen Grad statistisch gleich verhalten werden. Knotengruppen mit einem bestimmten Grad werden dann wieder mit der Meanfield-Methode untersucht. Diesen Analysemethoden ist gemeinsam, dass sie sich eigentlich nur das erwartete Verhalten der Prozesse ansehen. Und das erwartete Verhalten kann sich sehr stark von wirklichen Verhalten unterscheiden!

3 Der Voter-Prozess

Wenn man sich das wirkliche Verhalten von solchen Zufallsprozessen ansehen will, wird es sehr schnell schwierig. Ein Beispiel für die Probleme bei einer mathematischen Analyse von Markov-Kette liefert der sogenannte Voter-Prozess [5]. Er wurde eingeführt, um das Verhalten von Wählern (englisch: voter) zu studieren. Er ist lange bekannt und vergleichsweise einfach (Abb. 2). Jeder Knoten v wählt sich zufällig einen Nachbarn w aus und fragt ihn: „Was denkst du über einen bestimmten Sachverhalt?“ Mathematiker bilden Meinungen meistens durch Farben und manchmal auch durch Nummern ab. Hat w zum Beispiel eine Meinung, die in der Abbildung mit der Farbe Gelb symbolisiert wird, übernimmt v daraufhin die Meinung „Gelb“. Der erste Knoten adaptiert also ohne weiteres Nachdenken die Meinung des Nachbarn. Bei der Definition des Voter-Prozesses geht man davon aus, dass alle Knoten parallel die Meinung eines zufällig ausgewählten Nachbarn adaptieren.

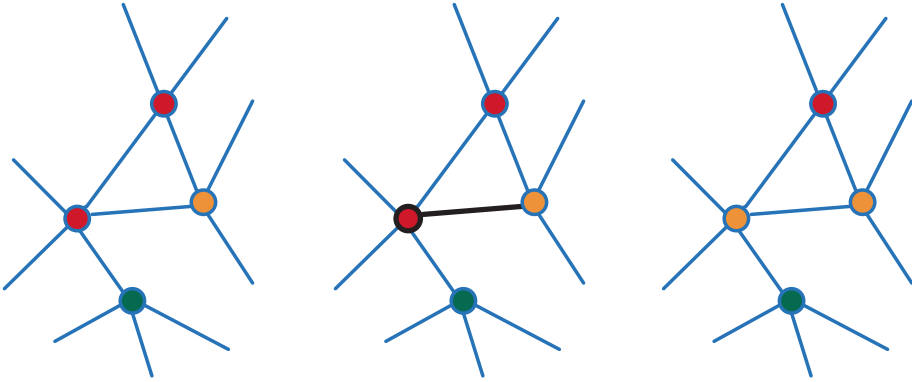


Abb. 2. Blick auf einen Knoten während des Voter-Prozesses: Er übernimmt die Meinung „Gelb“ eines Nachbarn

Der Voter-Prozess lässt sich einfach als Markov-Kette modellieren. Doch bei seiner Analyse gibt es zunächst zwei Schwierigkeiten. Erstens ist das Verhalten des Prozesse enorm abhängig vom zugrunde liegenden Graphen: Schon eine sehr kleine Änderung des Graphen oder des Prozesses selbst führt dazu, dass der Prozess komplett anders verläuft. Das macht eine mathematische Analyse der Prozesse sehr schwierig, weil sie diese Kleinigkeit mitmodellieren muss.

Zweitens steigt die Zahl der Zustände explosionsartig – mathematisch ausgedrückt: exponentiell – mit der Zahl der Knoten. Für 10 Knoten und drei Meinungen („Gelb“, „Rot“, „Grün“) ergeben sich $3^{10} = 59.049$ Zustände. Für 100 Knoten und drei Meinungen gibt es so um die 10^{48} Zustände. Zum Vergleich [6]: Das Gewicht der Sonne beträgt $2 \text{ mal } 10^{30}$ Kilogramm. Auf der Erde leben geschätzt 10^{30} Bakterien. Auf der Erde finden sich $6 \text{ mal } 10^{49}$ Atome. Das bedeutet: Die Zahl der Zustände einer Markov-Kette mit 100 Knoten und drei Meinungen beträgt ein 60-stel der Zahl der Atome auf der Erde. Der Faktor 60 ist angesichts der gewaltigen Zahlen zu vernachlässigen. Insofern hat die betrachtete Markov-Kette so viele Zustände, wie man Atome auf der Erde findet. Und doch ist die eine Markov-Kette mit 100 Knoten so klein, dass man mit ihr nicht viel darüber lernen kann, wie sich Meinungen oder Infektionen auf der Erde ausbreiten.

Für das nächste Problem, das sich bei der Analyse ergibt, ist es hilfreich, sich noch mal die Definition der Begriffe Erwartungswert und Varianz in Erinnerung zu rufen. Wir wollen uns das an einem sehr einfachen Beispiel ansehen: Zwei Menschen lassen einen Münzwurf darüber entscheiden, wer von ihnen die 1000 EUR auf dem Tisch bekommt. Bei Kopf gewinnt der eine Spieler, bei Zahl der andere. Der erwartete Gewinn eines jeden Spielers ist 500 EUR. Der echte Gewinn ist natürlich entweder 1000 EUR oder gar nichts. Der erwartete Gewinn ist genauso groß bei dem folgenden Spiel: mit Wahrscheinlichkeit $\frac{1}{2}$ erhält man 499 EUR und mit der gleichen Wahrscheinlichkeit 501 EUR. Die 500 EUR entsprechen jeweils dem mathematischen Erwartungswert, der nichts anderes ist als ein Mittelwert, gewichtet mit den Wahrscheinlichkeiten des Eintretens der einzelnen Ereignisse (hier Kopf oder Zahl). Beide Spiele haben also genau den gleichen erwarteten Gewinn. Aus Sicht

einer Mathematikerin liegt das daran, dass die sogenannte Varianz, definiert als die erwartete Abweichung vom Mittelwert, bei dem Experiment mit 1000 EUR sehr groß ist.

Bei unserem Voter-Prozess ändert sich der mathematische Erwartungswert nicht: Jeder Knoten übernimmt einfach die Meinung von einem zufällig ausgewählten Nachbarn – und das macht jeder Knoten zur gleichen Zeit. Erwartet ändert sich der Anteil der Knoten mit einer gewissen Meinung dadurch gar nicht. Doch tatsächlich ist das nicht so: Wenn man den Prozess beobachtet, haben nach kurzer Zeit plötzlich alle Knoten genau die gleiche Meinung. Danach passiert nichts mehr, kein Knoten kann seine Meinung ändern. Woran liegt das? Allein an der Varianz des Prozesses! Zunächst passiert, wie erwartet, erst mal nicht viel. Aber nach kurzer Zeit bekommt eine Meinung zufällig ein wenig mehr Unterstützung. Diese Meinung gewinnt dann eine kurze Zeit später.

Noch mal zurück zu unseren SIR-Prozessen. Bei den Analysen wird die Varianz meist vernachlässigt, man trifft also eine idealisierte Vorhersage. Oder auch eine Vorhersage, welche im Mittel eintreffen wird.

4 Der Majority-Prozess

Ein zweites Beispiel für die Schwierigkeiten der Modellierung epidemiologischer Vorgänge liefert der Majority-Prozess. Dieser ist dem Voter-Prozess ähnlich und immer noch sehr einfach. Bei ihm fragt jeder Knoten v (jede Person) nicht nur einen seiner Freunde, sondern zwei seiner Freunde (wieder repräsentiert durch benachbarte Knoten) nach ihrer Meinung. v hat nun seine eigene Meinung und er kennt die Meinung seiner zwei Freunde, also insgesamt drei Meinungen. Falls darunter zwei gleiche Meinungen sind, so übernimmt v diese sogenannte Majority-Meinung. Falls nicht, so gibt es zwei mögliche Varianten: Bei der ersten Variante bleibt v einfach bei seiner alten Meinung. Bei der zweiten Variante des Prozesses übernimmt er eine zufällige Meinung, beispielsweise die des ersten Freundes.

Jetzt stellt sich also die Frage: Verbreiten sich die Meinungen bei den zwei Varianten unterschiedlich? Berechnet man die jeweiligen Erwartungswerte, so lautet die Antwort: Nein. Die Prozesse würden sich demnach komplett identisch entwickeln. Doch tatsächlich verhalten sie sich völlig unterschiedlich. Bei der ersten Variante dauert es sehr lange, bis sich eine Meinung im ganzen Netzwerk durchgesetzt hat. Bei der zweiten Variante entwickelt sich der Prozess sehr schnell hin zu einem Zustand, bei dem im Netzwerk nur eine Meinung herrscht. Der Grund dafür ist die größere Varianz dieses Prozesses.

Wie aber werden Prozesse wie die Verbreitung von Infektionen, aber auch andere SIR-Prozesse in der Psychologie, der Physik oder der Informatik analysiert? Üblicherweise durch Differentialgleichungen. Bei Differentialgleichungen geht man implizit davon aus, dass man sich keinen Zufallsprozess mehr anschaut, sondern einen deterministischen Prozess. Anders ausgedrückt: Berechnet wird das erwartete Verhalten des Prozesses. Doch schon die Analyse des Voter-Prozesses hat gezeigt, dass diese Berechnung wenig aussagekräftig ist, vor allem dann, wenn die Infektionszahlen klein sind.

5 Soziale Netzwerke als Graphmodelle

Die Resultate, die wir in diesem Beitrag betrachtet haben, gelten meistens nur für reguläre Graphen, in denen alle Knoten den gleichen Grad haben. Soziale Netzwerke sehen demgegenüber weitaus komplizierter aus: Einige wenige Menschen haben sehr viele Kontakte, oftmals auch überall auf der Welt, da sie viel reisen. Daneben gibt es viele Menschen, die nur wenige und auch sehr lokale Kontakte haben. Ein weiteres Kennzeichen von sozialen Netzwerken ist eine Vielzahl an sogenannten Clustern und Dreiecken. Schließlich ist der maximale Abstand zwischen zwei Knoten, der sogenannte Durchmesser des Netzwerks, in sozialen Netzwerken sehr gering. Dies spiegelt sich in dem Ergebnis „Six Degrees of Separation“ (Sechs Grade der Trennung), die auch außerhalb der Informatik recht bekannt ist. Demnach sind alle Menschen auf der Erde nur sechs oder weniger soziale Verbindungen voneinander entfernt.

Das Modellieren von diesen sehr komplizierten sozialen Netzwerken ist ein Problem für sich. Eines der ersten theoretischen Modelle, das einige der oben genannten Charakteristika aufweist, wurde 1998 vorgestellt (Abb. 3). Die Knoten bilden einen Kreis, wobei jeder Knoten seinen Nachbarn im Kreis kennt. Einige Kanten gehen zu zufälligen Knoten quer durch den Kreis, wobei kurze Kanten wahrscheinlicher sind als lange Kanten. Dieses Modell ist ein recht guter Graph für Internetverbindungen und ist viel untersucht worden. Er kann aber auch als einfaches Modell für die Verbindungen zwischen Banken oder die Wechselwirkungen von Proteinen in den Zellen unseres Körpers dienen.

Watts-Strogatz model $N=20, K=4, \beta=0.2$

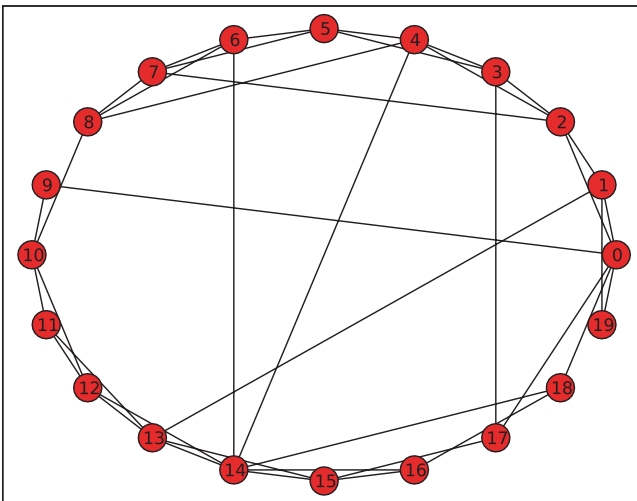


Abb. 3. Watts and Strogatz model (Arpad Horvath/Wikimedia Commons)

Wie kann man realistischere Modelle von Netzwerken finden? Das folgende Beispiel (Abb. 4) mögen Mathematiker sehr gerne. Es handelt sich um einen Graphen, der die Kooperation zwischen Mathematikern zeigt, in diesem Falle von Mathematikern mit der Erdős-Zahl 2: Paul Erdős ist ein sehr bekannter Mathematiker. Er selbst trägt die Erdős-Zahl 0. Jeder Forscher, der mit ihm eine wissenschaftliche Veröffentlichung geschrieben hat, trägt die Erdős-Zahl 1. Und jeder, der eine Publikation geschrieben hat mit jemandem, der eine Veröffentlichung mit Erdős geschrieben hat, erhält die Erdős-Zahl 2 zugeordnet. Dieser Kooperationsgraph ist typisch für ein soziales Netzwerk. Allerdings hat er 2100 Knoten und ist somit winzig gegenüber einem Netzwerk, auf dem man eine globale Epidemie analysieren oder modellieren sollte. Diese Modelle sind sehr kompliziert und sie lassen sich nicht einfach bilden. Dennoch hängt das Resultat oftmals sehr stark von dem benutzten Modell ab.

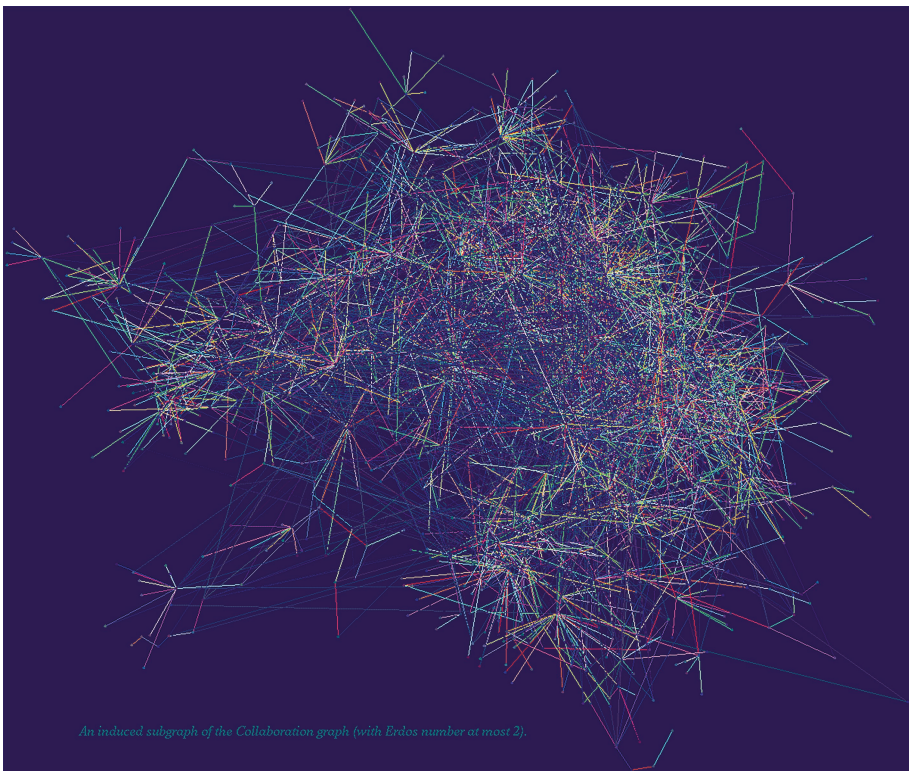


Abb. 4. Graph der Kooperation von Mathematikern mit der Erdős-Zahl 2 (Fan Chung Graham, <https://www.math.ucsd.edu/~fan/graphs/collaboration/>)

6 Simulation

Die vorangegangenen Erläuterungen geben einen Eindruck davon, wie schwierig die Analyse von epidemiologischen Prozessen ist. Leider ist deren Simulation auch nicht einfach. Die erste Hürde besteht darin, ein entsprechendes Netzwerk zu generieren. Selbstverständlich kann man ein Netzwerk entwickeln, indem man beispielsweise hundert Menschen fragt, wen sie kennen, um ihre Nachbarschaft zu bilden. Eigentlich ist man aber an sehr großen Netzwerken interessiert. Hat man auf diese Weise ein Netzwerk aufgebaut, so existiert es wahrscheinlich in dieser Form bereits nicht mehr, weil die Menschen andere soziale Kontakte gebildet haben oder ihre Stellung gewechselt haben. Auch aus diesem Grund benötigt man eigentlich ein theoretisches Modell, das SIR-Netzwerke generiert, um darauf dann die Simulation laufen lassen zu können. Informatiker arbeiten an diesem Problem schon sehr lange und entwickeln ein Modell nach dem anderen, aber das Ergebnis ist meist nur teilweise befriedigend.

Die zweite Hürde besteht darin, dass eine sehr kleine Änderung des Netzwerkes oder des Verbreitungsprozesses enorme, oft schwer nachvollziehbare Auswirkungen auf die Simulation haben kann. Oft ist es kaum möglich, den Prozess theoretisch ausreichend zu analysieren. Insbesondere werden häufig Varianzen nicht in die Analysen einbezogen. Wenn man aber einen Prozess nicht theoretisch analysieren kann, hat man ihn eigentlich auch nicht richtig verstanden. Dann sind auch die Resultate der Simulationen, die auf dem Netzwerk beruhen, nur äußerst schwer zu interpretieren.

Die dritte Hürde ist, dass die Graphen und die Zustandsmenge gigantisch groß sind. Dadurch sind die Simulationen sehr zeitaufwendig. Oft müssen sie einige Tage laufen – und erst dann merkt man womöglich, dass man das Ergebnis nicht versteht.

Literatur

1. Muchnik, L., et al.: Origins of power-law degree distribution in the heterogeneity of human activity in social networks. *Sci. Rep.* **3**, 1783 (2013)
2. Faloutsos, M., Faloutsos, P. Faloutsos, C: SIGCOMM '99: Proceedings of the Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication, S. 251–262 (1999)
3. Batagelj, V., Mrvar, A.: Some analyses of Erdős collaboration graph. *Soc. Netw.* **22**(2), 173–186 (2000)
4. Pastor-Satorras, R., Castellano, C., Van Mieghem, P., Vespignani, A: Epidemic processes in complex networks. *Rev. Mod. Phys.* **87**, 925 (2015)
5. Draief, M.: Epidemic processes on complex networks. *Phys. A Stat. Mech. Appl.* **363**(1), 120–131 (2006)
6. <https://www.physicsoftheuniverse.com/numbers.html>. Zugegriffen: 1. Febr 2021

Open Access Dieses Kapitel wird unter der Creative Commons Namensnennung – Nicht kommerziell – Keine Bearbeitung 4.0 International Lizenz (<http://creativecommons.org/licenses/by-nc-nd/4.0/deed.de>) veröffentlicht, welche die nicht-kommerzielle Nutzung, Vervielfältigung, Verbreitung und Wiedergabe in jeglichem Medium und Format erlaubt, sofern Sie den/die ursprünglichen Autor(en) und die Quelle ordnungsgemäß nennen, einen Link zur Creative Commons Lizenz beifügen und angeben, ob Änderungen vorgenommen wurden. Die Lizenz gibt Ihnen nicht das Recht, bearbeitete oder sonst wie umgestaltete Fassungen dieses Werkes zu verbreiten oder öffentlich wiederzugeben.

Die in diesem Kapitel enthaltenen Bilder und sonstiges Drittmaterial unterliegen ebenfalls der genannten Creative Commons Lizenz, sofern sich aus der Abbildungslegende nichts anderes ergibt. Sofern das betreffende Material nicht unter der genannten Creative Commons Lizenz steht und die betreffende Handlung nicht nach gesetzlichen Vorschriften erlaubt ist, ist auch für die oben aufgeführten nicht-kommerziellen Weiterverwendungen des Materials die Einwilligung des jeweiligen Rechteinhabers einzuholen.

