# Applications of complexity theory in official statistics

**2018 | 01**

Frank P. Pijpers

**Statistics Netherlands (SN) has a great interest in increasing its ability to offer more in the future than just tables of numbers. For instance SN is already regularly asked to offer expert analysis of causal connections between diverse statistical indicators. Social and economic phenomena are the reflection of complex interactions that occur within these domains, to a large extent out of sight of official statistics. Better knowledge of these interactions can make the difference between a correctly targeted pro-active and effective policy to solving a social or economic problem, or a reactive policy. The issue with the latter is that often only after resources have already been used (or wasted) can it be determined whether that commitment has been useful. Based purely on chance, if there is insufficient targeting of a policy the likelihood of effective use of resources will be low.**

**The central SN task of reporting on the 'state of the country' must be done in a much more refined way, certainly also as targets to improve sustainability rise higher on the social and governatorial agendas. The aim of complexity studies within the framework of national statistical institutes such as Statistics Netherlands is among others to be able to indicate correctly which sectors of business activity are vulnerable/ fragile and also which are most robust. The same aims arise in a social setting: to determine which social groups are isolated or vulnerable. Another aspect is the behaviour in time of such a system: if there are sudden large changes in the system, what early warning indicators can be designed that signal such changes? The standard statistical output is not sufficient for this goal. The toolkit of complexity studies can give substance to this strategic goal, and to the social added value of national statistics and Statistics Netherlands in the longer term.**

**One immediate result of the exploration of applications of complexity theory by SN is the realisation that the form of the distribution function for the turnover of Dutch businesses is very stable over the four years 2010-2013. It can be characterized with only a few parameters, which are quite different from the way these data usually are presented.**

# 1 Introduction

Statistics Netherlands (SN) has a long tradition in descriptive statistics. The institute is best known for its publication of official statistics in the form of tables, with short explanatory messages for the media. In addition to the mandatory statistics program, it is expected that SN will focus more on offering a package of customized services to local and central governments. For openness of public administration and for fact-checking by the media and the public, publications of custom-made tables and open data initiatives remain important. However, it must be acknowledged that exclusively descriptive statistics offer insufficient added value and that there is also a need for more in-depth analyses. Many important clients of SN will not themselves have the necessary expertise available, while SN has sufficient expertise to offer appropriate and extensive technical support.

Social phenomena and developments are the result of interdependent underlying mechanisms and processes. It is precisely these mechanisms that can often be uncovered with the aid of statistics, where there is a need for the governments and other parties that SN can and should want to serve. It is important to choose carefully the steps in this statistical analysis process, tailored to the problem, and to be able to accurately describe and justify these choices when

reporting the results of that process. Although SN certainly already has a lot of expertise in this, it is essential to gain further knowledge of and experience with new domains of methodology.

Important methodological developments lie in the field indicated by the term 'complexity' or 'complex systems'. This branch of science is concerned with the behavior of systems that consist of a multitude of individual parts (or 'actors') that influence each other. In (Buiten et al., 2017) an overview can be found of specific complexity approaches that appear to be relevant to particular issues that SN needs to tackle as national statistical institute. The way in which research in this area is tackled is necessarily model-based, where data is used in statistical testing of the assumptions with which the model is constructed. The word *model* is not used here in the sense of an ideal image, nor is it a composite of images and behaviors of each individual component that are as detailed as possible. The term model is meant here as a simplified description of the system or the parts inside it, with a minimum number of fixed / chosen properties. By parametrising the model features the researcher has means to control, simulate and adjust the behavior of this model system. The model must also produce measurable outcomes: it must be possible to observe it, albeit virtually. The next requirement is that on these virtual observations of the model, it must be possible to apply statistical measurement methods, in exactly the same way as the real economy or real society can be observed and statistically analyzed. In this way part of the behavior of reality can be reproduced approximately. The usual situation is that there is also a part of the behaviour of the real system that can not be reproduced. Such failures then provide guidance to improvements on the model. Its use thus lies in incrementally uncovering and describing the 'real' mechanisms.

Statistical tests are an essential part of the modeling technique. Because a test can indicate where the model is statistically significantly different from reality and must therefore be rejected, it can also be deduced on which points a model *is too simple* and needs to be expanded. In this sense, *technical* exploration of the consequences of (government) policy within model-based research is equivalent to *a well-founded and quantitative statistically tested statement* about the validity range of a simplified model: such phased model building as knowledge is built up is commonplace in for instance the natural sciences.

While model-based estimation is already being used within the working practices of SN, this usually still takes place in a restricted meaning. In some cases, the model refers to a specific stochastic process, such as a Poisson process or a normally distributed process. In other instances, the properties of a previous observation period, or of a similar but different (geographical) domain are used as a model for a new period or a domain that is not yet observed. At the cross section between complexity theory and official statistics, the model has a much more central role than this, because the research question that is asked is more focused on, for example, how and especially why a system has a certain reaction after an (external) stimulus.

Examples of applications of complexity theory can be found in both the behavior of the Dutch economy, as well as of wider society. In the context of the 'sustainable development goals' (SDG) cf. (Smits and Hoekstra, 2011) there are important questions that can be directly linked to complexity. How should it be assessed, for example, whether a certain economic activity, or set of related activities, is sustainable? Within complex systems there is a known phenomenon referred to by the term 'self-organized criticality' (SOC) (see chapter 9 of Downey (2012) for an introduction). It appears that many systems have the property that, through the interactions of the elements within the system, they reach a critical state. As long as there is no external stimulus or disruption to this system, it appears to be in a stable equilibrium. However, a very small disturbance of the system in this critical state can cause very far-reaching changes in the

system. In this case a traditional approach, where the behaviour over the recent past is used as an indicator of future stability or robustness, is clearly not suitable. Some small stimuli might cause widespread disruption, and others none at all. An arbitrary external stimulus of a system that is in such a critical condition will not as a matter of course produce an instability or other 'catastrophic' behavior. Modeling is a means of assessing in what ways, ie. for which incentives or stimuli, an economy is 'vulnerable/fragile' and which stimuli are harmless, i.e. leave the economy 'robust' or 'sustainable'.

An example of a link between complexity and broader social phenomena lies in the areas of social cohesion and aging. The extent to which people are or are not isolated within society can be a big influence in the extent to which they have to rely on professional or social organizations. The extent to which individuals or groups feel connected to social life, or feel represented at and by public institutions, can also be driven by the branching of their network of social interactions. Network theory is an important methodological framework that is relevant to quantify how much or how little cohesion social networks have, and can therefore also be an instrument to be able to identify where and how fast social cohesion changes.

There are also issues where economic and social problems come together, such as the determination of the influence of (economic) inequality on the combined economic activity of a population. An article that looks at the self-limiting behavior of growth in an economy from a complexity perspective, which shows that a greater degree of inequality can be a disruptive influence on growth is eg Jerico et al. (2016).

With these two examples, two important methodological branches within complexity theory are also mentioned. In the following sections, each of these two is explained in more detail, namely *dynamic system theory* on the one hand, and *network theory* on the other. A good introduction that describes the analysis and modeling of complex systems by Downey (2012). A recent inventory of currents and representative research articles and books is Newman (2011). The specific needs of SN are addressed in greater detail in (Buiten et al., 2017).

# 2 Dynamic system theory

## 2.1 frame of reference

In dynamic system theory a complex system is modeled as consisting of entities with only a few properties. Each of these elements can interact with any other element, where the result of that interaction may be that the two elements that have had the interaction have received a different value for those properties. In fact, this is a translation from the natural sciences where the properties of a material are statistically determined because the material consists of a large number of molecules: it usually concerns numbers typically of the order of $10^{24}$. In a gas or liquid molecules can move freely and therefore collide with other molecules: they have an interaction. In a collision energy and momentum are exchanged: the energy and momentum properties of each molecule after the collision have other values than they had before.

In the context of economics, this approach is often described as 'agent-based' modeling. The role of molecules is here taken over by the 'agents'. Those 'agents' can be companies, or consumers,

or workers. When, for example, companies are these agents, the interaction can be an exchange of services or goods and money in business-to-business trade. With such an interaction, the property 'turnover' of both partners gets a larger value, and possibly also a property such as 'profit' may gain a larger or smaller value. Such an approach has for example been used to model the character of the growth of companies in the US: see eg Ormerod (2002). One important difference compared to the behaviour of molecules is that the number of companies per country is probably of the order of $10^5$ or so, and even globally it is not more than $10^8$. These numbers are sufficiently many orders of magnitude smaller than the $10^{24}$ molecules considered before, that it makes a genuine difference in terms of modelling the collective behaviour.

With this method of modeling, it is tempting to continue the analogy with molecules even further, although there are (economic) grounds to remain cautious about this. In the simplest modeling of an 'agent based' / dynamic system, it is assumed that the 'agents' have no memory. That means the outcome of — or the relative probabilities of the different possible outcomes of — an interaction is completely independent from previous interactions that the agents have had. Only the current state of the agents (ie the vector of values of the set of properties that agents in the model have) that are involved in an interaction, then plays a role in the calculation of the interaction result. This approach offers certain advantages in the ease with which models can be simulated, but at the same time it is a strong limitation. For example, companies will not in practice do business with arbitrary trading partners. Instead they will normally enter into contracts and have fixed trade agreements for some length of time with a fairly stable selection of partners. Translated to a model, this means the opportunity for an agent to interact with another agent will be larger if there has been an interaction with the same agent before. It is also possible that when in the past a conflict has arisen between companies, a new contract is likely to be avoided. Also in this situation there is therefore a need for 'memory' among the agents in a model, because the chance of an interaction in this latter case is greatly reduced. As indicated in Gallegati et al. (2008), it is important that existing economic knowledge is explicitly embedded in models. In the context of SN, this means the next steps in this research will be taken in close collaboration between methodologists and various specialist economics departments.

This does not alter the fact that simple modeling can be used to assess to what extent 'memory' plays a measurable role in the dynamics of the Dutch economy. By way of illustration of these principles, a simple 'toy model' is presented.

## 2.2  toy model

To illustrate a description of the company population in the Netherlands as a dynamic system, a very simple model is constructed. The goal in this context is not primarily to get a true view of the company population, but to show which steps are needed in building a method:

1. Characterizing the agents (companies), ie. the answer to the question of what characteristics of a company in economic connection are most relevant to their trade activity. In a realistic model, properties such as a sector of main activity and any secondary activities, numbers of personnel, and working capital undoubtedly also play a role.
2. Formulating the interaction between companies in the form of rules cq. the dependencies of the probability of an interaction to the characteristics mentioned under point 1. Likewise the result of such an interaction in terms of state variables must be formulated. In the context of SN these interactions together are referred to as the intermediate consumption.

3. Quantifying the influence of an external world: ie. the consumer market, international competition, and legal provisions and restrictions. Within SN, the term final spending is used when it comes to, among other things, consumption, export, inventory mutations, and investments by companies. The added value of complexity research is particularly relevant here.

In the design of this toy model, every company is seen as an agent with as its only characteristic property the *turnover*. Companies can therefore only differ from one another in the size of their turnover, and the most important outcome of the model is the distribution function $\psi$ of the turnover for the business population. This is a genuinely measurable quantity, since SN has a number of relevant registers at its disposal in which, among other things, an annual turnover is registered per company per year. In what follows the word business will be avoided as much as possible, and the word agent will be used, explicitly to continue to indicate that it concerns properties, relationships, and behavior as adopted within the model context. It is recognized that at the moment the model is still far away from the 'real life' company population.

The interaction that agents have in the model is an exchange of goods and / or services, each of which can be expressed in terms of a certain monetary value $x$ (goods) resp. $y$ (services) which together constitute the turnover. The distribution function $\psi$ can thus be formulated as a function of this $x$ and $y$: $\psi(x, y)$. This can also be formulated as a total turnover $\rho$ with a relative distribution within it among goods and services such that:

$$x \equiv \rho \cos \phi$$
$$y \equiv \rho \sin \phi \qquad (1)$$

where the angle $\phi$ is a measure of the relative distribution of the turnover between traded services and goods. With this choice the distribution function becomes a function of $\rho$ and $\phi$: $\psi(\rho, \phi)$. Formally this is here defined as:
*$\psi(\rho, \psi)$ is the aggregate turnover of agents with individual turnovers between $\rho$ and $\rho + d\rho$ and relative distribution between goods and services within it, which is mathematically processed as an 'angle', between $\phi$ and $\phi + d\phi$.*

The next step is to specify the interaction between agents. As described above, the most simple modeling is obtained if it is assumed that the interaction between agents only depends on the turnover size $\rho$ itself. For example, it can be assumed that the chance for two agents to have an interaction is strongly peaked around a turnover difference of $0$ between the agents. That means that:

- agents prefer to act with other agents with a comparable turnover size $\rho$, so that after a (positive) interaction both may have a modest increase in turnover size, ie. modest as a percentage of their turnover.
- agents have mainly competition from other agents with a comparable sales size $\rho$, so that one (negative) interaction means that one or even both may have a modest decrease in turnover size.

This type of interaction, largely limited to agents with almost equal $\rho$, resulting in only modest changes in $\rho$ per interaction, translate to a diffusion equation for the turnover distribution function. For an isolated system in equilibrium, without consumers or an external market, then the differential equation for $\psi$ applies:

$$\nabla \cdot (D \nabla \psi) = 0 \qquad (2)$$

where the diffusion coefficient coefficient $D$ depends on $\rho$ and $\psi$. A large value of $D$ means that turnover can be redistributed quickly over the total population of agents, while a small value means that the redistribution is slow. For simplicity here it is further assumed that D will only depend on $\rho$. A next step is a system that is not perfectly balanced, but that can have a slow overall growth or shrinkage of turnover with time. In that case a term must be added to the equation (2) with the time derivative $\partial \psi / \partial t$:

$$\frac{\partial \psi}{\partial t} = \nabla \cdot (D \nabla \psi) \tag{3}$$

For simplicity, it is customary to assume that the time dependency can be expressed as $\partial \psi / \partial t = \alpha \psi$.

If the system is not isolated, it means that for example in (3) there is a source term on the righthand side of the equal sign in cf. (3). For example, a form $S(\rho)$ can be used: the influence of the external market is such that for agents with a turnover between $\rho$ and $\rho + d\rho$ agents merge or disappear so that the combined turnover increases at a rate that is given as a 'source' term $S$. If $S > 0$ the external factors stimulate the activity or the continued existence of agents, while as $S < 0$ business activity is suppressed. In what follows, also for $S$ it is assumed that it depends only on $\rho$ and not on $\phi$, just as is done for $D$. Writing out all terms in the differential equation produces:
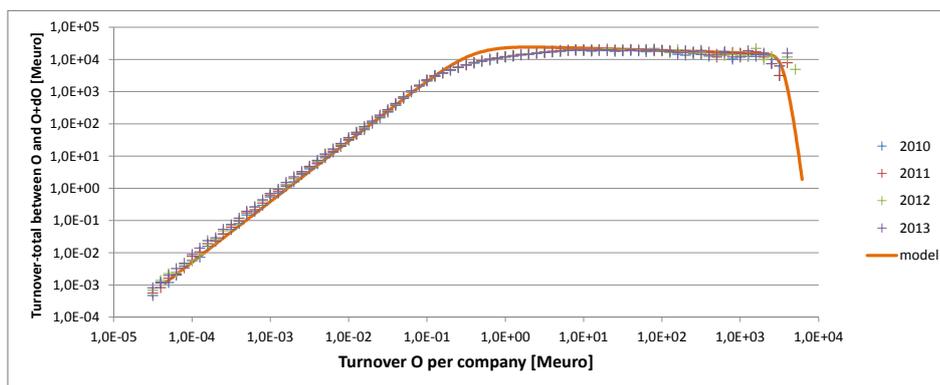
$$\alpha \psi = \frac{1}{\rho} \frac{\partial}{\partial \rho} \left( D(\rho) \rho \frac{\partial \psi}{\partial \rho} \right) + \frac{D(\rho)}{\rho^2} \frac{\partial^2 \psi}{\partial \phi^2} + S(\rho) \tag{4}$$

If it is further assumed that the main component of $\psi$ is also independent of $\phi$ the relevant term disappears, so that the differential equation can be written as:

$$\rho \frac{\partial}{\partial \rho} \left( D(\rho) \rho \frac{\partial \psi}{\partial \rho} \right) = \rho^2 \left[ \alpha \psi - S(\rho) \right] \tag{5}$$

By giving a prescription for $D(\rho)$ and $S(\rho)$ and then solving the differential equation, the model distribution function $\psi(\rho)$ is determined. Conversely, with a chosen $D$ and a measured $\psi$ it is also possible to check which external conditions (source function) $S$ leads to that form for $\psi$.

**Figure 2.1 The distribution function $\psi$ for the turnover, the model (solid line) as well as the data from the companies register, for the years 2010 through 2013**



Within SN, diverse business data are available on a yearly basis from a combination of registers and business surveys. For the purpose of this discussion data are used for the four calendar years 2010 through 2013. A selection has been made in the sense that companies that are in fact self-employed people without personnel (ZZP), present as businesses in the register, are removed from the population. The turnover of companies is then placed in a histogram, but with much finer bins than is usual for standard presentation by SN on eg Statline. Such a histogram thus contains the numbers of companies $N(\rho)$ with a turnover between $\rho$ and $\rho + d\rho$. The

distribution function can be determined from: $\psi \equiv \rho N(\rho)$. For each of the years, $\psi$ is shown in fig. 2.1. Also included in the figure is a fit function that is a power law over the range of 0 to 0.2 M€, and a separate different power law across the range of 1 to $2 . 10^3$ M€, ie. it has the form $\psi(\rho) \propto \rho^\gamma$ over each range with different values for $\gamma$, with a smooth transition between these two ranges. Above a company turnover of $2 . 10^3$ M€ there exist only a few companies so not every bin of the histogram is still filled, and this is modeled by a steep drop of $\psi$. For the low turnover range, the value of $\gamma$ seems to be 1.95, while for the higher range $\gamma = -0.073$. In the transition area between the two ranges the fitting function (in orange) is clearly too 'sharp' compared to the actual $\psi$. Overall, this type of shape which is a combination of power laws is characteristic of known systems. Natural phenomena where "self-organized criticality" of the system plays a decisive role for distribution functions very often display such power law behaviour. For the US population of companies such a size distribution was also examined, on the basis of a sample: see eg Axtell (2001).

What is clear from this figure 2.1 is that this form is very stable over the four years, and can be characterized with only a few parameters: the slopes in the two ranges (ie the two values of $\gamma$), the place of the transition between those ranges, and the 'cut-off' point. Those are not parameters that SN would traditionally choose to present, but perhaps they are the parameters that would be of most interest for external researchers.

The analysis of the equation (5) does not stop at this point. It is useful for the further steps to define a variable $\zeta$:

$$d\zeta \equiv \frac{1}{\rho D} d\rho \tag{6}$$

and substituting in Eq (5) so that it can be written as:

$$\frac{\partial^2}{\partial \zeta^2}\psi = \rho^2 D \left[\alpha\psi - S(\rho)\right] \tag{7}$$

where $\rho$ can be implicitly determined from $\zeta$

$$\zeta = \int \frac{1}{\rho D} d\rho \tag{8}$$

By using the definition of the Fourier transform $\Psi$ of the distribution function $\psi$ (Bracewell, 1965):

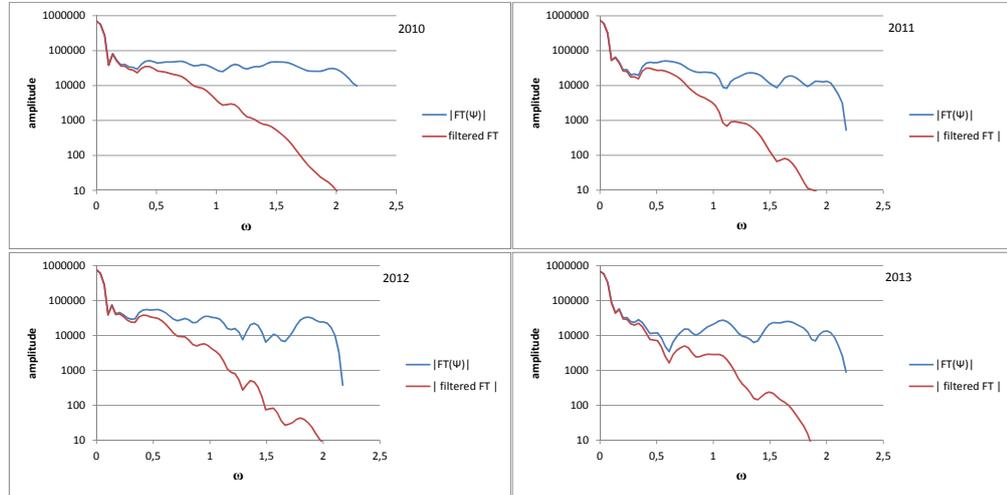$$\Psi(\omega) \equiv \frac{1}{\sqrt{2\pi}} \int \psi(\zeta) e^{i\omega\zeta} d\zeta \tag{9}$$

it can be shown that:

$$S(\rho) = \alpha\psi - \frac{1}{\rho^2 D} \frac{1}{\sqrt{2\pi}} \int -\omega^2 \Psi(\omega) e^{-i\omega\zeta} d\omega \tag{10}$$
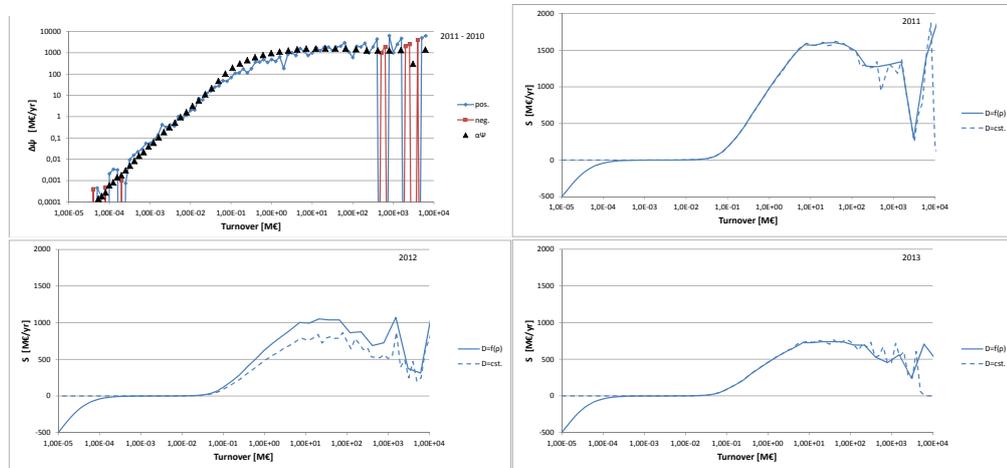
where the integration in Eq. (10) is the inverse Fourier transform. Numerical methods for the execution of the (discrete) Fourier transformations (DFT) in Eq. (9) and (10) are trivial. Note that this analysis does not constitute solving a differential equation. The approach (10) is taken because $\psi$ is measured only at discrete intervals, and therefore its second derivative is poorly defined. By taking the DFT, multiplying by $\omega^2$ and taking the inverse, a second derivative is obtained numerically and this difficulty is circumvented. Because of the multiplication of $\Psi$ with $\omega^2$ under the integral sign in (10), the determination of $S$ becomes very sensitive to the value of $\Psi$ at large values of $\omega$ and thus is very sensitive to irregularities cq. a stochastic component in $\psi$. This is precisely the same as would happen when using a simple finite difference method for approximating a second derivative. It is necessary to apply some filtering between the Fourier transformation in (9) and the inverse in (10) to suppress the absolute value of $\Psi(\omega)$ for the

largest values of $\omega$. For instance in the top lefthand panel of fig. 2.3 the influence of the rapid variations at the high end of company turnovers above a few hundred M€ is suppressed. The effect of how the filter suppresses the Fourier transform at high $\omega$ is shown in fig. 2.2.

**Figure 2.2    Panel top left: The unfiltered (blue) and filtered (red) Fourier transforms of $\psi$ for the years 2010 up to and including 2013. $\omega$ is in arbitrary units.**
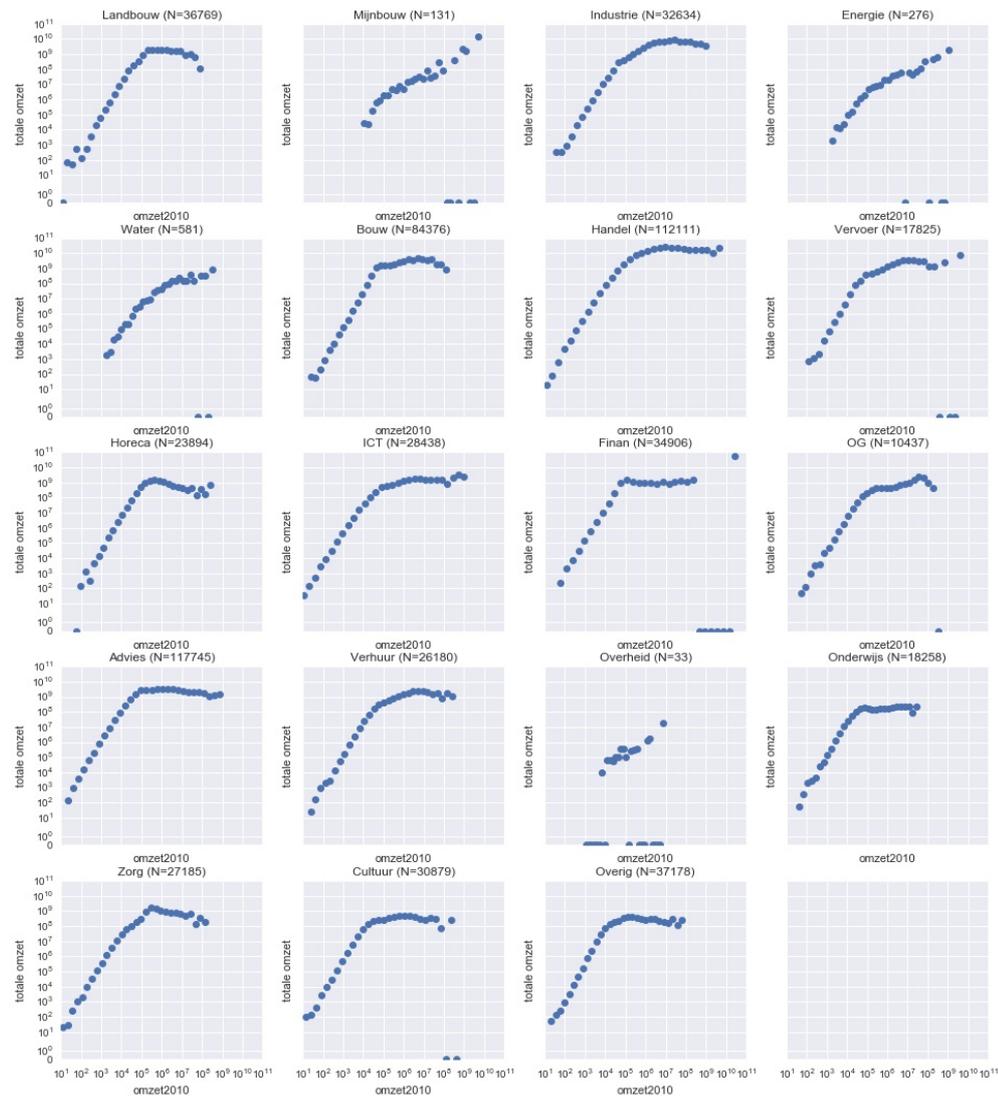


**Figure 2.3    Panel top left: The growth $\psi(t) - \psi(t-1)$ between 2010 and 2011 (positive in blue, negative in red), with the fit $\alpha\psi$. Other panels: the source functions $S$ for the years 2011 up to and including 2013, for two different choices for $D$. When the source function $> 0$ the external factors favour growth, whereas when the source function $< 0$ the external factors inhibit growth.**



From equation (10) it appears that it is possible to determine from the observed distribution function $\psi$ what the source function $S$ is, ie. for what range in business turnover the market is stimulating or inhibiting growth. The parameter $\alpha$ can be determined for each year by a least-squares fit of the year-on-year difference $\psi(t) - \psi(t-1)$ with $\psi(t)$ itself. The result of such an analysis is shown in figure 2.3, for two different assumptions for $D$. The simplest assumption is that $D(\rho)$ is constant. The results shown are for a value $10^4$ €$^2$/yr. Various different values between $10^3$ and $10^8$ show slight differences but are not visible at this scale. A second possibility that is straightforward to explore, is a linear relationship $D = D_0 \left[1 + b \ln(\rho/\rho_0)\right]$. Here results are shown for $D_0 = 10^4$ €$^2$/yr, $\rho_0 = 0.1$€ and $b = 1$. When using this behaviour for $D$ the results are more sensitive to the value of $D_0$ than if $D = D_0$ for all $\rho$, in particular at low values of $\rho$. Larger values of $D_0$ will make $S$ turn downward to large negative values at higher values of $\rho$.

The parameter $\alpha$ is determined for the consecutive years, 2010-2011, 2011-2012, 2012-2013, and is resp. $\alpha = 7.9\,\%$ growth, $\alpha = 5.1\,\%$ growth, $\alpha = 3.8\,\%$ growth. The source functions also show that for companies with turnover above 50 k€ in these years the market is clearly less stimulating, but still positive, but that for companies with turnovers below this limit, the source function over the course of all of these years is 0 or negative.

**Figure 2.4    The distribution function $\psi$ for the turnover for various Dutch business sectors, for 2010. Top to bottom, left to right: agriculture, mining, (heavy) industry, energy, water, construction, trade, transportation, accomodation and food services, ICT, finance, real estate, consultancy, rental & lease, public administration, education, health care, culture, other**



An interesting extension of this is to separate out the various business sectors to see whether differences can be seen between sectors. Fig. 2.4 shows the distribution function $\psi$ for the year 2010 (J. Roest, internship report). From this it is clear that in many sectors the distribution function has a similar shape as $\psi$ foor all sectors combined, although in detail the slopes over the subranges and the location and 'sharpness' at the turnover point differ between the sectors. There are a few exceptions to this pattern: mining (top row, second panel in 2.4), energy (top row, last panel), water (second row, first panel), and government administration (fourth row, third panel). These four business sectors may be fundamentally different from the others, but they are also much smaller in the numbers of companies per sector. Further research is still

ongoing at SN and is the subject of future reporting.

## 2.3 broader interest

A 'toy' model as discussed in the previous section has obvious shortcomings, since assumptions are made that are evidently not an accurate reflection of reality. It does have value to illustrate how the idea of an agent based economy can be translated into a mathematical model with which quantitative model results can be produced. The growth rates $\alpha$ give a representative, ie. well-weighted, image of turnover growth for the entire population of companies which is a key indicator for the state of the Dutch economy. The further analysis of the source function can be used to assess, even per business sector, at which size of companies a stimulus package or support measures are most desirable in times of economic 'heavy weather'. It also serves as an evaluation means to assess afterwards whether such a package has had the intended effect.

Comparable analyses can be performed on the income distributions and asset distributions of all residents of the Netherlands. In this case, this could provide a clearer picture of the relative tax burden for different population groups. Furthermore, any discussions about a balanced and equitable distribution of that relative tax burden can be given a better factual foundation.

# 3 Network theory

In the previous section, the dynamic system theory approach was presented, in which the actors or agents in the system have no 'memory', unless explicitly taken into account in the modeling. When modeling from the point of view of network theory an opposite starting point is used. The role of the agents in system theory is here taken over by the nodes in a network. If there are connection(s) between two nodes, this indicates that between those nodes there can be interactions. Nodes that are not directly connected in this case never have a direct interaction, although there may still be indirect influences if they do belong to the same connected network.

In the simplest, static, modeling of companies in the Dutch economy in the form of a network, there is no room for a company to be able to interact (conclude a contract) with another company with which there has not been an interaction before. When dynamics have to be taken into account, this means explicitly steps in modeling must be included, in which new connections can be placed in the network, and also existing connections can be broken. A paper that explicitly considers dynamical networks in the context of economic activity and the most recent debt crisis is (Dehmamy et al., 2014).

As with dynamic systems, the nodes in the network are interconnected, and the network itself is actually a mapping of the route(s) along which nodes exert influence on each other. A phenomenon that is specifically suited to the network analysis framework is to provide a description and modeling phenomenon of *cascades*: a process that reinforces itself so that a relatively small disturbance very quickly causes very large changes over large parts of a network (Motter and Yang, 2017). Examples are a localized traffic accident that causes many tens of kilometers of traffic jams, or the failure of an electricity distribution station, resulting in disruption of the electricity supply of a large city and thus also a large part of public transport, but also a social media video that goes 'viral'.

## 3.1 The topology of networks

The first step of useful applications of network theory is 'simply' mapping connections. From this map or image follow conclusions about the structure of the network, also called the *topology*. In most cases not every node has a connection with all other nodes in the network. There are various (statistical) measures for how many connections each node has to other nodes. If there is a large number of connections for the majority of nodes, it will mean that the *redundancy* in the network is high. Depending on the phenomenon described by the network, that can that be a favorable or an unfavorable characteristic. For example, when it comes to high voltage cables that distribute electricity over the country, at least some redundancy is desirable. If any particular cable or link fails due to a technical defect, that redundancy ensures that the energy supply to a particular node can be maintained via another route through the network. However, when it comes to a network of infections, for an infectious disease, a large network redundancy implies that through the large number of connections between nodes a disease can quickly take on epidemic forms. In some cases it is necessary to take direction into account in the connections between nodes ("node A affects node B but not vice versa"). The structural characteristics of a network, whether this is a social or economic network, are clearly directly relevant to quantifying vulnerability, isolation, integration, and influencing. See (Buiten et al., 2017) for further detail.

Another type of network are hierarchical networks or modular networks. In both types the distribution of numbers of connections per node is skewed: there are relatively few nodes with many connections, while there are many nodes with just one one or a few connections. Some key features of such hierarchical networks are described in articles about hierarchical networks (Gao et al., 2010; Benson et al., 2016). A specific kind of hierarchical network is a tree structure, such as a folder structure in a computer system, or a set of decision rules used for classification or in a data cleaning process.

It must be noted that mapping out all connections in a network can be difficult or even impossible, especially as the number of nodes increases. For this reason, sometimes a modeling approach is used where the missing information is compensated for by assuming certain characteristics for the distribution of connections. One way to then determine the most likely configuration is by defining an entropy for network configurations and maximizing that entropy. It can be shown that great care must be taken in defining the entropy however. Two ways in which this is typically done using concepts from statistical physics, can be shown to lead to different behaviour of the entropy. In physics this phenomenon occurs when particles have many long range interactions. In other words, nodes in a network can still have a strong influence on each other, even if their minimal separation is a number of links away in the network. Therefore, conversely, if one node has behaviour that appears strongly correlated to another, one cannot automatically assume that there must be an (unobserved) direct link between the nodes. Further details can be found in (den Hollander, 2016) and (Garlaschelli et al., 2016).

## 3.2 a toy network model

To illustrate the relevance of network theory, it is useful to explore again by means of a toy model, in which ways it can produce results for national statistical institutes. It is clear from section 3.1 that the structure of the network itself is already interesting. Interactions that occur in the network are also of interest when calculating the time evolution of the system, for example to be able to determine whether, and what kind of equilibrium situation exists in the

system. In this framework, the interactions between companies over a certain period of time are described with the aid of a transition matrix $G$. If $X$ is a column vector with a status variable per company (employee numbers, turnover or another business property relevant for official statistics), this can change by interactions in the network so that at the transition from time $t$ to time $t+1$ the vector of state variables changes from $X_t$ to $X_{t+1}$. Under the assumption that this transition only depends linearly on the state variable this is described as:

$$X_{t+1} = G \cdot X_t \tag{11}$$

The $N \times N$ transition matrix $G$ is then a reflection of the network topology: if between the nodes $k$ and $k'$ in the network there is no direct connection, its elements are $G_{kk'} = G_{k'k} = 0$. In the sense that this transition between the states at times $t$ and $t+1$ only depends on the state at time $t$ and not on previous states of the system, this is a Markov process. If the state variable is a business turnover, and the interaction involves exchanging money, goods and services, the transition matrix $G$ perhaps has to have certain symmetry properties. If for example, with a particular transaction the exchange of money, goods and services between companies $i$ and $j$ is exactly equal in value, so neither of the two companies involved makes a profit, then:

$$G_{ij}X_{j\,t} = G_{ji}X_{i\,t} \tag{12}$$

The time evolution of such a system then consists of repeatedly applying the matrix $G$ to the state vector $X$, making it relevant and interesting to determine which properties this matrix $G$ has.

Note that there is a conceptual difference between this matrix $G$ and the (Leontief) input-output table / matrix, in which the vector $X$ would be the total production per sector instead of turnover per individual company, and $G \cdot X$ would then be the intermediate consumption. That interpretation would require in addition a term for final output (final demand), a vector $d$, on the righthand side of eq. (11). However, the form used here is closer to a concept of a circular economy, where private households are simply another sector: purchasing goods and services from the other sectors, but also producing recyclable materials or energy for re-use in any (other) sector. In this way, the vector $X$ simply gains an element representing the sector 'private households', and the matrix $G$ has an extra row and column. One might then even include (finite) natural resources as an element in the vector $X$ where the corresponding extra row in $G$ has all elements $= 0$, except the diagonal element which is between 0 and 1.

For this discussion it is useful to be able to separate out $G$ into a symmetrical and an anti-symmetrical component:

$$
\begin{aligned}
G &\equiv S + A \\
S &\equiv \tfrac{1}{2}\left[G + G^T\right] \\
A &\equiv \tfrac{1}{2}\left[G - G^T\right]
\end{aligned}
\tag{13}
$$

where the superscript $T$ means taking the transpose of a matrix. By using these definitions for the matrices $A$ and $S$, $S$ is a symmetric matrix $S^T = S$ and $A$ anti-symmetrical $A^T = -A$. Symmetrical matrices can always be factorized as follows:

$$S = U\,D\,U^T \tag{14}$$

where the matrix $D$ is diagonal, and $U$ is a unitary matrix: all its columns are mutually perpendicular and have vector length 1, so that

$$U^T U = I \tag{15}$$

The diagonal elements of $D$ are also called the *singular values* of the matrix $S$, and the column vectors of $U$ the associated singular vectors. It is customary to arrange the elements of $D$ and

therefore also the vectors $U$ so that the absolute values of the diagonal elements of $D$, the $|d_{ii}|$, run from large to small with increasing $i$. Often, this decrease is monotonous, but in principle there is the possibility that the system is degenerate, ie. that there are singular values that are identical to each other.

In a system where $G$ itself is already symmetrical, and therefore $G = S$ and $A = 0$, it is easy to see from combining eqs. (11), (14), and (15) that:

$$X_{t+n} = G^n \cdot X_t = UD^nU^TX_t \tag{16}$$

This is important because $D^n$ like $D$ itself is diagonal, with ordered diagonal elements $d_{ii}^n$. That means that as $n$ gets bigger, the diagonal elements decrease faster and faster with increasing index $i$ on the diagonal: the ratio with respect to the first element $|d_{ii}/d_{11}|^n \downarrow 0$ for increasing $n$. After a while such a process evolves the system to a situation where there is a fixed ratio between the elements in the vector of state variables $X$. That fixed ratio is described by the first of the singular vectors $U_1$, unless the largest singular value is degenerate. That is the stable equilibrium ratio. The overall size of all the elements of $X$ grows or shrinks all together in the course of time, depending on whether $|d_{11}| > 1$ or $< 1$. Every initial deviation from this stable equilibrium will disappear exponentially with time. The characteristic time scale depends on the projection of that deviation on the other singular vectors. The the most slowly disappearing deviation has a characteristic time scale $\tau = -1/\ln|d_{22}/d_{11}|$ as long as $d_{22} \neq d_{11}$. If at microscopic level the relationships between companies evolve towards an equilibrium, this will also apply to a distribution function. as shown in fig. 2.3.

A direct consequence of this analysis is that when a certain economic phenomenon, of which the time evolution can be described by a Markov process, does not seem to have a stable equilibrium but, for example, cycles or a different kind of temporal evolution, then this must be the result either of the presence of an anti-symmetrical component in the transition matrix $G$, or it is the case that two or more of the largest eigenvalues are equal to each other: $d_{11} = d_{22} = \ldots$.

If the matrix $G$ is purely anti-symmetric, i.e. $G = A$ and $S = 0$ then it follows that:

$$\left(G^2\right)^T = (AA)^T = A^TA^T = (-A)(-A) = A^2 = G^2 \tag{17}$$

This means that almost the same analysis steps follow as described above for a $G$ which is purely symmetrical, but everywhere where $G$ itself is listed, $G^2$ must be used. In the event that $G$ is anti-symmetric, it can be proven that all singular values are equal to $0$ or are purely imaginary. For a $N \times N$ anti-symmetric matrix $A$ where $N$ is odd, there is at least one singular value $= 0$, or an odd number. All other singular values are pairs of complex conjugates, with real parts $= 0$.

In the general case when $G$ has both a symmetrical and an anti-symmetrical contribution, the singular values are complex ie. with both real and imaginary parts. When $N$ is odd, there is at least 1 purely real singular value, and generally an odd number, while the remainding singular values are pairs of complex conjugates. When $N$ is even, there are $0$ or an even number of real singular values, and the rest consists of pairs of complex conjugates. Incidentally, the actual singular values can be degenerate: the same real value can occur several times, in this case the simple inference as followed above is no longer valid. There is no balance for the network, where the system strives for a fixed ratio of turnover described by the singular vector, belonging to the largest (in absolute value) singular value. It can be seen from the purely anti-symmetrical example the whole system in that case has a period 2 because $G^2$ is symmetrical, and therefore

only has real singular values. For the general matrix with a spectrum of complex singular values, there are multiple periods in the system. The matrix $G$ can still be factorized:

$$G = U \, D \, U^T \tag{18}$$

but now $D$ is diagonal with a combination of $m$ real singular values $\mu$ and $K = (N - m)/2$ pairs of complex conjugate singular values $\lambda e^{\pm i\theta}$:

$$D = \begin{pmatrix}
\mu_1 & 0 & .. & .. & .. & .. & .. & .. & .. & .. & .. \\
0 & \mu_2 & 0 & .. & .. & .. & .. & .. & .. & .. & .. \\
.. & 0 & .. & 0 & .. & .. & .. & .. & .. & .. & .. \\
.. & .. & 0 & \mu_m & 0 & .. & .. & .. & .. & .. & .. \\
.. & .. & .. & 0 & \lambda_1 e^{i\theta_1} & 0 & .. & .. & .. & .. & .. \\
.. & .. & .. & .. & 0 & \lambda_1 e^{-i\theta_1} & 0 & .. & .. & .. & .. \\
.. & .. & .. & .. & .. & 0 & \lambda_2 e^{i\theta_2} & 0 & .. & .. & .. \\
.. & .. & .. & .. & .. & .. & 0 & \lambda_2 e^{-i\theta_2} & 0 & .. & .. \\
.. & .. & .. & .. & .. & .. & .. & 0 & .. & 0 & .. \\
.. & .. & .. & .. & .. & .. & .. & .. & 0 & \lambda_K e^{i\theta_K} & 0 \\
.. & .. & .. & .. & .. & .. & .. & .. & .. & 0 & \lambda_K e^{-i\theta_K}
\end{pmatrix} \tag{19}$$

The ranking on the diagonal is now done by magnitude in absolute value, separately for the $\mu$ and the $\lambda$. For example, it is quite possible that $\lambda_1$ is greater than all $\mu$. Where in an anti-symmetrical matrix the complex phase $\theta$ is always equal to $\frac{\pi}{2}$, that is not the case now, and every value in the open interval $(0, \pi)$ occur. Each complex phase $\theta_i$ is associated with a periodicity with period $2\pi/\theta_i$ of the system, ie. the matrix $G$ must be multiplied by itself this number of times, for that particular singular value to become real. When this happens the associated column vector in $U$ is reproduced, multiplied by the real factor $\lambda_i^{2\pi/\theta_i}$. Since all $\theta_i$ may be different, it is a priori unlikely that there exists an $\alpha$ such that $D^\alpha$ has only real values on the diagonal.

Following the same reasoning as above for a symmetric matrix $G$, it can be concluded that if $\mu_1 > \lambda_1$ all cyclical phenomena will eventually be suppressed. If, however, $\mu_1 \leq \lambda_1, ..., \lambda_j$ for a certain value of $j \leq K$ there is no static equilibrium, but cycles persist in the business turnover system. There is then no fixed ratio where the system converges, but it continues to exhibit oscillations, possibly with various periods. The dominant periodicity can be identified with a business cycle, which therefore does not occur because of limitations imposed externally to the system, but is an emergent phenomenon of the interactions between the companies in the population.

## 3.3 broader interest

Also in the case of the toy model from the previous section there are limitations in how realistic the modeling is. For example, it is assumed that the transition matrix $G$ does not change over time, and also that there is no dependency between time $t + 1$ and times further back in time than $t$. This has, among other things, the direct consequence that the only systems that can be modeled must have exponential growth or shrinkage (boom or bust), or show oscillations.

In a realistic system, the transition rate $G$ itself will also change over time. At least there ought to be provision for a stochastic contribution, so that a system will never fully converge. In combination with eigenvalues that are degenerate, or almost degenerate, such a disturbance term can easily cause a system to switch in a stochastic manner between two different states described by different singular vectors. This is an example of a system that is at least to the eye in

a state of low-dimensional chaos. If there are non-linear effects that play a role in limiting the (relative) growth of companies that are part of a network, they can also play a role. Due to nonlinear effects, the matrix $G$ can change slowly over time, in such a way that some or even several singular values come close together in the course of time. Again this leads to a system that is in a state of low-dimensional chaos.

For official statistics, the identification of business cycles, and the place of the Dutch economy in a cycle, is an important indicator of great governmental and economic interest. If there is sufficient data available about relevant interactions between companies, an analysis is possible with which balances or periodicities can be measured. In addition this framework can be used to map the dynamic sensitivity to disruptions of the network of companies. This means there is a direct relevance for quantifying a concept such as a vulnerability in the economy.

If this concept is used instead to model the (circular) economy at a mesoscopic level, eg per business sector, then the interpretation of the largest eigenvalue becomes of interest. Such a system has losses: energy loss and waste that is not recycled and hence value that is being destroyed. Presumably this has the consequence that the largest eigenvalue is $< 1$ in absolute value, given finite natural resources. The largest eigenvalue of $G$ could be an indicator of sustainability of the economy. Conversely, if a requirement of stability, and hence a value of $1$ for the largest eigenvalue is imposed, this can also be interpreted as a way to assign economic value to the remaining natural resources.

# 4  Conclusion

This discussion paper indicates that two important paradigms within complexity theory, dynamic system theory on the one hand and network theory on the other hand, are directly relevant to qualitative concepts such as sustainability, and economic equilibrium. Modelling the real world using these concepts can help to quantify positioning the Dutch economy in the business cycle and to measure vulnerability / fragility. When this knowledge is applied to the topology of networks of people and their social interactions, it becomes possible to quantify integration between social groups as well as isolation of minorities or other social subgroups.

The current main output of SN is publishing tables of numbers, for example on Statline. This evidently remains necessary, so that also for external parties it remains possible to apply their own analyses to the best available data. There is, however, a difference between these numbers themselves and their information content. SN can only meet the needs of external parties if the information content locked up in the numbers is exposed. SN can not afford to continue to use only the traditional tools that have been used for this purpose, but must actively pursue extensions, such as those from complexity theory, to its arsenal of output.

An important clear example is that the form of the distribution function for the turnover of Dutch businesses is very stable over the four years 2010-2013, and can be characterized with only a few parameters: the slopes of this function in two subranges of turnover size, the place of the transition between those ranges, and the 'cut-off' point. Those are not parameters that SN would traditionally choose to present, but perhaps they are the parameters that would be of most interest for external researchers.

A second example, from transition matrices in a network, shows that the largest eigenvalue of such a transition matrix could be an indicator of sustainability of the (circular) economy.

Conversely, if a requirement of stability, and hence a value of 1 for the largest eigenvalue were to be imposed, this can also be interpreted as a way to assign economic value to the remaining (finite) natural resources.

Although qualitative objectives will always have a use and relevance within administrative considerations, increasingly government policies benefit from targeted and accurate measures and quantitative indicators. This is certainly the case, among other things, for the sustainable development goals. Better techniques are essential to be able to assess whether measurable and structural progress is being made in achieving those objectives. The means, methods and techniques that are available, based on complexity theory, will play a central role and are indispensable in a package of output and services from national statistical institutes.

# References

Axtell, R. (2001). Zipf distribution of US firm sizes. *Science 293*, 1818--1819.

Benson, A., D. Gleich, and J. Leskovec (2016). Higher-order organization of complex networks. *Science 353*, 163--166.

Bracewell, R. (1965). *The Fourier transform and its applications*. McGraw-Hill. (paperback edition: 1999).

Buiten, G., E. de Jonge, and F. Pijpers (2017). Het cbs en complexiteitstheorie: een position paper. Technical report, CBS.

Dehmamy, N., S. Buldyrev, S. Havlin, H. Stanley, and I. Vodenska (2014). Classical mechanics of economic networks. *arXiv preprint arXiv:1410.0104*.

den Hollander, F. (2016). Netwerken bekeken vanuit de statistische fysica. *Nederlands tijdschrift voor natuurkunde november*, 372--375.

Downey, A. (2012). *Think Complexity*. Green Tea Press. (second edition).

Gallegati, M., S. Keen, T. Lux, and P. Ormerod (2008). Worrying trends in econophysics. *Physica A: statistical mechanics and its applications 370*, 1--6.

Gao, J., S. Buldyrev, S. Havlin, and H. Stanley (2010). Robustness of a network of networks. *arXiv preprint arXiv:1010.5829*.

Garlaschelli, D., F. den Hollander, and A. Roccaverde (2016). Ensemble nonequivalence in random graphs with modular structure. *Journal of Physics A: Mathematical and Theoretical 50(1)*.

Jerico, J., F. Landes, M. Marsili, I. Castillo, and V. Volpati (2016). When does inequality freeze an economy. *arXiv preprint arXiv:1602.07300*.

Motter, A. and Y. Yang (2017). The unfolding and control of network cascades. *Physics Today 70*, 32--39.

Newman, M. (2011). Complex systems: A survey. *Am. J. Phys. 79*, 800--810.

Ormerod, P. (2002). The US business cycle: power law scaling for interacting units with complex internal structure. *Physica A: statistical mechanics and its applications 314*, 774--785.

Smits, J. and R. Hoekstra (2011). Measuring sustainable development and societal progress: overview and conceptual approach. Technical report, CBS.