

UNIVERSITY OF PADOVA

LECTURE NOTES
OF
LIFE DATA EPIDEMIOLOGY

COLLECTION OF THE LECTURES NOTES OF PROFESSORS CHIARA POLETTO AND SANDRO MELONI.

Author:
ALICE PAGANO

Revised by:
ANDREA NICOLAI



Thursday 12th November, 2020

Contents

I Meloni's Lectures	1
1 Basics Definitions and Compartmental Models	3
1.1 Comportamental models	3
1.2 Basic models	6
1.2.1 SI model	7
1.2.2 SIS model	8
1.2.3 SIR model	10
1.3 Extensions of the SIR model	12
1.3.1 SIR with Demography	12
1.3.2 SIRS Model	13
1.3.3 SEIR Model	14
1.4 Summary of compartmental models in well-mixed populations	15
2 Network Science - Basics	17
2.1 Main definitions	17
2.2 Degree distribution over networks	19
2.2.1 Erdős and Rényi Model: random graphs	19
2.2.2 Scale-free networks	21
2.2.3 Barabási-Albert Model	23
3 Epidemic Spreading on Networks	25
3.1 SIS model in a network	25
3.1.1 Homogeneous Networks	26
3.1.2 Heterogeneous Networks	28
3.2 SIR model in a network	31
3.2.1 Degree-based mean-field theories (DBMF)	31
3.2.2 Individual-based mean-field theories (IBMF)	32
3.2.3 DBMF vs IBMF: Epidemic threshold	34
3.2.4 IBMF and Pair approximation	36
4 Epidemic spreading on networks: advanced models	37
4.1 Markovian Models	37
4.2 Non-Markovian Epidemic Spreading	38
4.2.1 SIR Model with Multiple Infectious Stages	39
4.2.2 Generalized SIS Model	40
5 Spreading in social systems	45
5.1 Complex contagion	45
5.2 Applications to Online Social Networks	48
II Poletto's Lectures	51
Bibliography	53

Part I

Meloni's Lectures

1

Basics Definitions and Compartmental Models

Some models are wrong, but most of them are useful.

– Unknown author

Models in science have two different roles: **understanding** what happens and **predict** will happen. Models can be of two types: simple and more complex ones. In the simplest ones we just consider the minimal number of parameters and events involved: this indeed allows to understand what are the main mechanisms of a phenomenon.

In this course, we are going to start with very simple models in which we assume that there is no structure behind in the population. Obviously this is not accurate, but allows us to understand at a first glance some underlying mechanisms. Then, we are going to consider social structures and introduce contact network models. We will also take into account interactions among different populations and exploit data to understand how members move from one population to another. Finally, we are going to introduce deal with the so called “Agent Based” models, for a quick overview on them.

1.1 Comportamental models

We now introduce the **comportamental models**. These are fundamental since the most of epidemiological theories are based on them. In reality, however, there are different levels of understanding how diseases can diffuse: we can consider the disease only at a biological level, or at simpler one. Note as it is practically impossible to insert all the details of a process in a single model. We therefore need to summarize all the biological processes in few **parameters** which describe, on average, what we can see inside the population. This is the same principle behind the statistical mechanics in which we look for large scale (macroscopical) effects.

Let us consider a population of individuals and try to characterize it. Note as we have not made any assumption on the individuals and relationships between them. We now introduce three different **compartments**, denoted with **S** (that stands for *Susceptible*), **I** (*Infected*), **R** (*Recovered*), and want to label people according to the stage of their disease, as seen in Fig. 1.1. However, one should note that there can be also transitions from one state to another one, according to some rates that describe the **dynamic**. For instance, in Fig. 1.1 these are β and μ .

This approximation, on the other hand, is quite strong: by keeping the rates fixed we are assuming that the process underlying the spreading of the disease is a **Markovian**. In reality, we do not see exponential distributions, but a decay. **previ-**

Lecture 2.
Friday 2nd
October, 2020.
Compiled:
Thursday 12th
November, 2020.

ous sentence not clear This last point will be discussed during the course when we will deal with “**non-Markovian**” epidemics. The interpretation we may give to β is the “*per contact*” *infectious rate*, in this way we only need to count the number of contacts. Different models can be introduced according to the type of the disease: for instance **SI**, **SIR**, **SIS**, **SEIR** and so forth.

One should note that medical status is actually different from infectious status. In the latter we do not care about medical status of the person, but only about the disease and how the immune system reacts against it.

As an example, for the **SEIR** compartmental model, we have four main stages of the disease: starting from a healthy state (**Susceptible**), the individual can contract the disease (**Exposed**) and then, only after some time, becomes infectious (**Infectious**) until he recovers (**Recovered**) (Fig. 1.2). The most important thing to keep in mind is that these compartments are not the same ones of the medical status, since they keep into account different parameters despite the disease is the same one.

Now, let us introduce the **Basic Reproductive Number R_0** (pr. “*R naught*”) which is a measure of the infection in the population. If we wanted to empirically determine it: we put one guy inside a group for an arbitrarily long time period and, at the end, we count the number of secondary cases that we have. This is the main idea behind the computation of R_0 . By definition R_0 is the average number of secondary cases an average primary case produces in a **totally susceptible population**. This parameters therefore determines whether a disease will spread or not:

$$\begin{cases} R_0 < 1 \\ R_0 = 1 \\ R_0 > 1 \end{cases} \quad (1.1)$$

Let us consider the plot of Fig. 1.3, we have a sort of **second order phase transition** at the point $R_0 = 1$. Note that R_0 for the SARS is higher than the one of COVID-19. However, we did not experienced an outbreak of this disease, so that is not the only parameters to be taken into account in the models. In order to compute R_0 we assume that the population is totally susceptible. This is however valid only at the very early stages, later on, we must consider both epidemiological and demographical aspects. The conclusion is the following: R_0 may vary from one population to another.

Since we are doing a **coarse-graining** of the dynamics, this number represents the average of all possible different distributions. A wrong argument is to think that similar R_0 lead to similar outbreaks. The distribution of infections can be quite heterogeneous: the mean could be quite representative only if we are dealing with homogeneous populations, that is not the case for real networks. For instance, let us consider the plot in Fig. 1.4. We see that SARS is heterogeneous, while Spanish Flu was more a homogeneous one. COVID-19 is most likely somewhere in the middle.

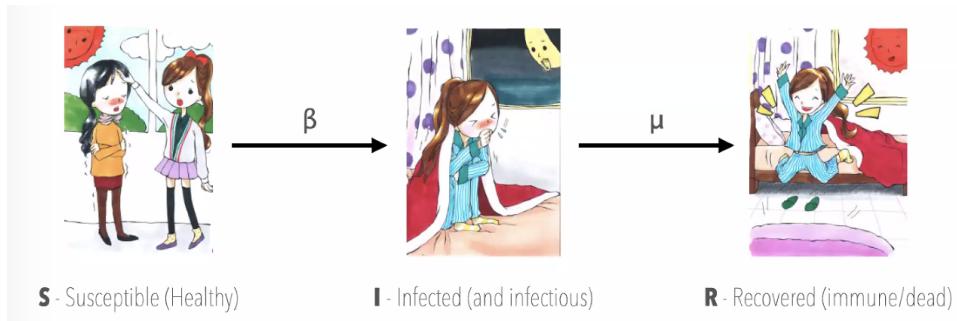


Figure 1.1: Classification of infected population in three different stages of the disease.

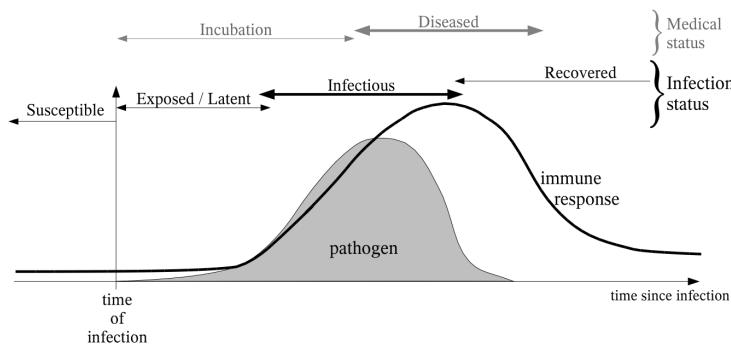


Figure 1.2: A sketch of the time-line of infection, showing the dynamics of the pathogen (grey area) and the host immune response (black line) with the labeling for the various infection classes: **Susceptible**, **Exposed**, **Infectious**, and **Recovered**. Note that the period when symptoms are experienced (medical status) is not necessarily correlated with any particular class of epidemiological models.

Let us now introduce the **Effective Reproductive Number $R(t)$** , which is the same of R_0 but varying wrt time. Hence, it is the average number of secondary cases that a single case produces in a population at time t .

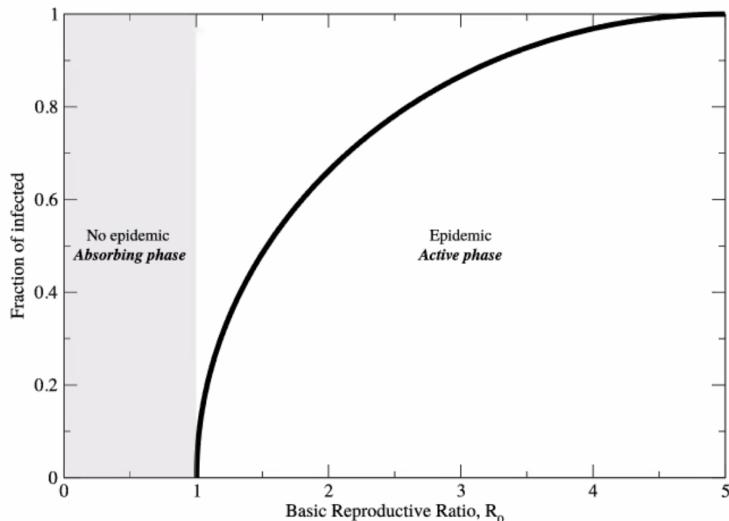


Figure 1.3: Fraction of infected vs basic Reproductive Ratio, R_0 .

Other important quantities we may want to introduce are:

- **Infectious period:** average period for a person to be infectious and is computed as $\tau = \frac{1}{\mu}$, $\tau = \frac{1}{(\alpha+\mu)}$ where the presence of α depends on the model (α : average duration of "exposed time" stay);
- **Incubation period:** period of time between infection to occurrence of symptoms;
- **Generation time:** time for an infected person to generate a second infection;
- **Serial interval:** time between the onset of symptoms for a person and the onset of symptoms for another second infected person;
- **TOST:** time between the onset of symptoms to an infection.

A problem in predicting a possible outbreak of a disease is that TOST in many cases can be negative (see Fig. 1.5 for more details). See Fig. 1.6.

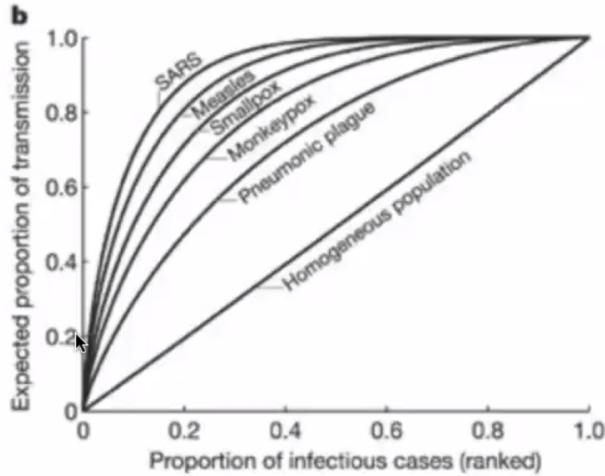


Figure 1.4: Figure from: Lloyd-Smith et al. Nature 438, 355–359 (2005).

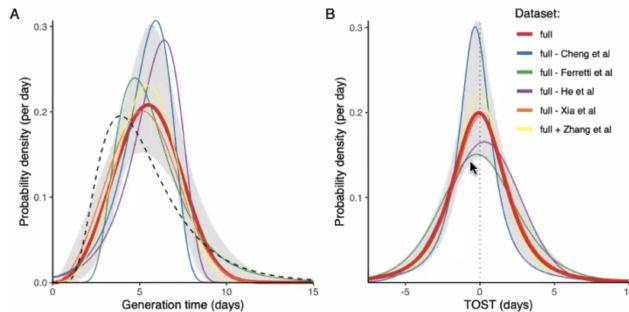


Figure 1.5: Figure from: Ferretti et al. <https://www.medrxiv.org/content/10.1101/2020.09.04.20188516v1>

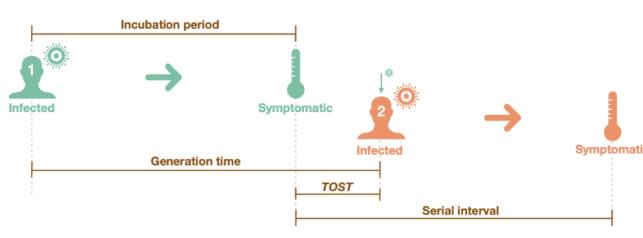


Figure 1.6: Figure which explains better the already defined quantities.

1.2 Basic models

In this lecture we are going to introduce some of the basic models we will use for the entire course. The first assumption we make is that we are in **well-mixed populations**, or in other words *homogeneous mixing*. Mathematically, it is what is called **mean field approximation**.

In the well-mixed population assumptions, it holds that that:

- all individuals are **equivalent**, hence every one has the same probability of being infected;
- every individual has the **same number of contacts** $N - 1$, or on average $\langle k \rangle$;
- We are in a **closed population**. That is to say that the sum of the density distribution of the individuals is equal to 1, hence we have no deaths or births.

In practice, we are assuming that our time scale is so little that we can consider the population constant.

1.2.1 SI model

The simplest model one can think of is the **SI** (**Susceptible Infected**). In this model one can get the infection and, once we have got, we cannot recover, that is to say we stay infected forever.

The **transition diagram** that describes this model is the following:



where β is the “*per contact*” *infection rate* and dictates the speed of the spreading. We can write down the **equation** that can be solved exactly:

$$\begin{aligned} \frac{ds}{dt} &= -\beta \langle k \rangle si \\ \frac{di}{dt} &= \beta \langle k \rangle si \end{aligned} \quad (1.3)$$

where $\langle k \rangle$ represents the average contacts, while i stands for the fraction of infected people in the entire population ($i = I/N$), and s is the fraction of susceptible people in the population ($s = S/N$). The product si is the probability of having a contact between an infected and a susceptible, and βsi is the probability of having a contact between an infected and a susceptible which in turns leads to an infection.

One of the most important quantity we may want to introduce in our lexicon is the so called *prevalence* $i = \frac{I}{N}$, that is another way to define the density of infected people wrt the entire population.

In order to solve it analytically, we recall that our population is closed. Therefore $s + i = 1$, and it follows that we only have one equation to be solved since $s = 1 - i$. We have that:

$$\frac{di}{dt} = \beta i(1 - i) \rightarrow \frac{1}{\beta i(1 - i)} di = dt \rightarrow \frac{1}{\beta(1 - i)} di + \frac{1}{\beta i} di = dt$$

Integrating both sides:

$$-\log|1 - i| + \log|i| = \beta(t + C) \rightarrow \frac{i}{1 - i} = e^{\beta(t+C)} = Ae^{\beta t}$$

with $A = i_0/(1 - i_0)$. The result is:

$$i(t) = \frac{i_0 e^{\beta t}}{1 - i_0 + i_0 e^{\beta t}} \quad (1.4)$$

which is a sigmoid function (Fig. 1.7) that always saturates at 1. One should note that after the first part, where the growth is actually exponential¹, then at a certain point the slope starts to decrease. The reason for this is that the contribution given by the term si , namely the probability of funding new susceptible people, decrease. Finally, we saturates at 1 after some. As can be clearly seen from Fig. 1.7, it is the value of β that drives the spreading. By increasing it, we obtain a faster exponential growth. This actually was the simplest model one can think of.

Remark. In the course we are going to use capital letter for integer numbers, while small letters refer to densities.

¹It is the one we have seen in the media for COVID-19.

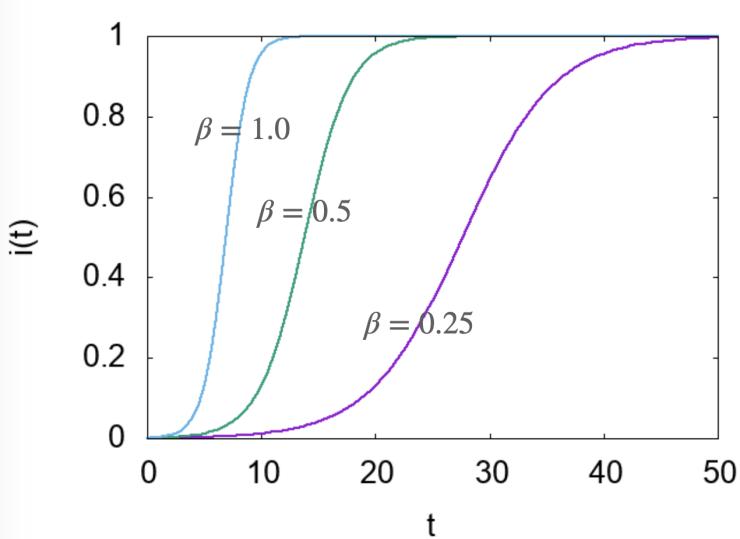


Figure 1.7: Plot of the solution of the SI model for different β .

1.2.2 SIS model

Now, let us introduce a slightly more complicated model, that is the **SIS** model, where compartments are **Susceptible**, **Infected**, **Susceptible**. Transitions now are two:



where the first transition is mediated by I , that is to say we need to encounter another infected to contract the disease, while the second one occurs **spontaneously** according to the rate μ .

This model is used for diseases that do not confer immunity. When we use the expression **endemic state** it means that the disease keeps on circulating in the population for very large times.

The most important feature about this model is that it is the simplest one where **dynamical equilibrium** can be reached. Therefore an individual may recover from the disease, but he does not get immunity. Indeed there are always people infected that can propagate the disease. The μ is the **recovery rate** which determines the *time-scale of the infection*. Dividing β by μ you can **rescale** all the **dynamics**. The **equations** are exactly the same as before, except for a term:

$$\begin{aligned} \frac{ds}{dt} &= -\beta \langle k \rangle si + \mu i \\ \frac{di}{dt} &= \beta \langle k \rangle si - \mu i \end{aligned} \quad (1.6)$$

and in addition can solved in the very same way we previously did.

Also, the shape of the **solution** is a sigmoid as before:

$$i(t) = i_0 \frac{(\beta - \mu)e^{(\beta-\mu)t}}{\beta - \mu + \beta i_0 e^{(\beta-\mu)t}} \quad (1.7)$$

By plotting it, one should note that despite the same form, we do not saturate at 1, but at $\frac{\beta-\mu}{\beta}$. Hence, as we said, we have some sort of **dynamical equilibrium**: the number of new infected is more or less the same of the new recovered people at each moment. The density $i(t)$ will therefore fluctuate around this value $\frac{\beta-\mu}{\beta}$ and, by enlarging μ , we can obtain larger fluctuations (Fig. 1.8).

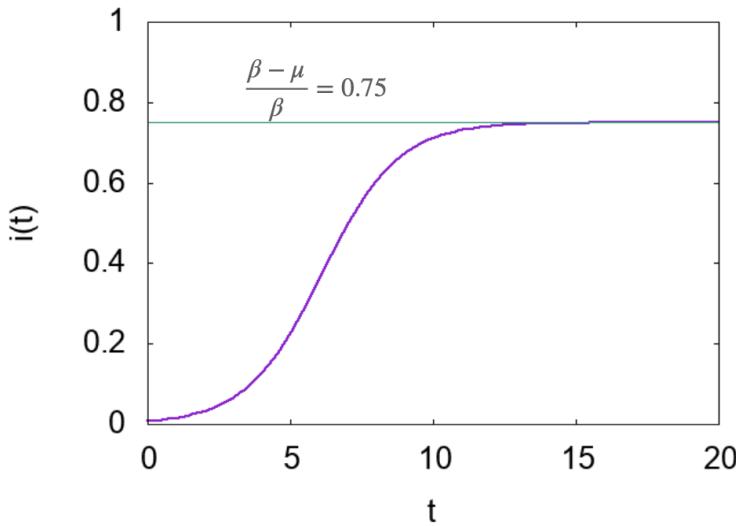


Figure 1.8: Plot of the solution of the SIS model.

It can be instructive to study what happens according to this model at the **transient**. At the beginning, one can assume that almost the entire population is composed by susceptible people ($s \sim 1$), while the number of infected is very small ($i \ll 1$). Hence, the differential equations can be rewritten as following:

$$\frac{di}{dt} = \beta \langle k \rangle si - \mu i \sim \beta \langle k \rangle i - \mu i \rightarrow i(t) \sim i_0 e^{(\beta \langle k \rangle - \mu)t}$$

One should note that if $\beta \langle k \rangle < \mu$ there is no spreading at this point anymore, while, if $\beta \langle k \rangle > \mu$ the exponent becomes positive and from this follows the exponential growth at the beginning (Fig. 1.9).

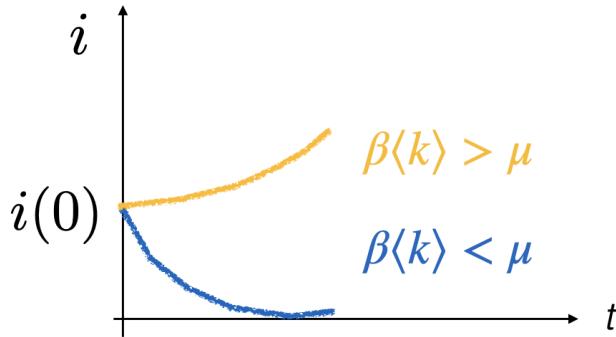


Figure 1.9: Initial transient for the SIS model.

One very important thing is that considering the **steady state** we can have two possible behaviors:

$$\frac{di}{dt} = 0 \rightarrow \begin{cases} i = 0 & \beta \langle k \rangle < \mu \\ i > 0 & \beta \langle k \rangle > \mu \end{cases}$$

and we have that:

$$i > 0 \iff \beta > \beta_c = \frac{\mu}{\langle k \rangle} \quad (1.8)$$

where β_c is known as the **epidemic threshold**. This tells us whether the disease is going to spread.

In addition the epidemic threshold is the minimum value of the infection probability for which the disease survives. This is what in physics is called a **second order**

phase transition (Fig. 1.10). In this case the **critical exponents** are the same of the Ising model, since they belong to the same class of universality. β_c is one of the most important quantities we are going to study.

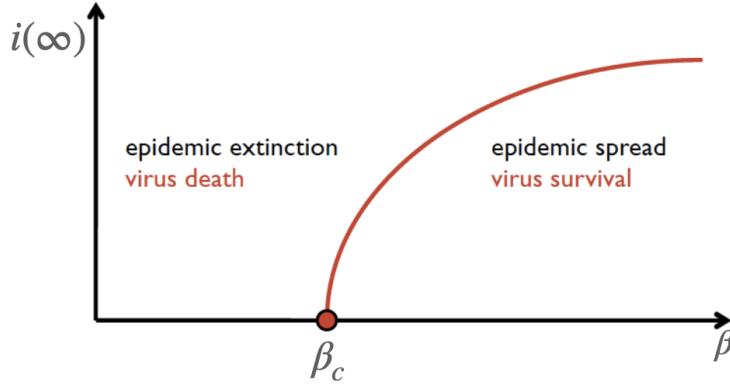


Figure 1.10: Epidemic diagram.

One may ask what is the relation between R_0 and the epidemic threshold. Obviously, they are strongly correlated. We actually say that given a **critical value**, below it we have no spreading, while above we have a fraction of infected people.

The epidemic threshold is given by the condition under which we observe the spreading. Mathematically, given a specific model, its critical version will return the values of the parameters for which $R_0 = 1$. If we are slightly above this threshold, we need a minimum of infected people and the disease is going to spread. Considering for instance the case of the SIS model:

$$R_0 = \frac{\beta}{\beta_c} = \frac{\beta \langle k \rangle}{\mu} = 1 \quad (1.9)$$

In general, for an epidemic to occur in a susceptible population R_0 must be > 1 , so the number of cases is increasing.

1.2.3 SIR model

We now discuss the so called *SIR* model, whose compartments are **Susceptible**, **Infected** and **Recovered**. The idea behind is the same one of the SIS, but we are now adding a new state which accounts for long lasting immunity (**R**). Hence, once a person has got the disease and has recovered, he obtains a long **immunity**. Recall that, since we assumed that the population is closed, its density is still fixed to 1.

The transitions for this model are:



and one should note that we cannot have any endemic state. For large times all individuals will have been infected, and recovered, so the disease will be spreading no more.

The differential equations that describe this model are:

$$\begin{aligned} \frac{ds}{dt} &= -\beta \langle k \rangle s i \\ \frac{di}{dt} &= \overbrace{\beta \langle k \rangle s i}^{\text{New infections}} - \overbrace{\mu i}^{\text{Recovery}} \\ \frac{dr}{dt} &= \mu i \end{aligned} \quad (1.11)$$

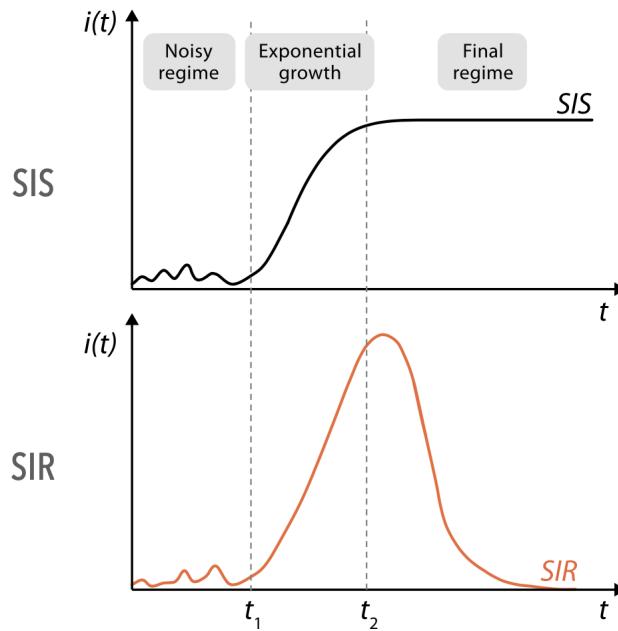


Figure 1.11: Epidemic regimes.

This is actually a good point to introduce the **different regimes** we may encounter during a spreading, which are represented in Fig. 1.11 for the SIS and SIR models.

Initially, at the beginning of each spreading, we see the so called **noisy phase** where the numbers are too small to cause a large spreading. Here we can observe only some sort of stochastic fluctuations. In many cases, we can end up without any spreading: if we are able to stop the infected people in this noisy phase, we stop the disease (assuming it is heterogeneous **not clear**). If it is not the case, the disease starts spreading according to the characteristic **exponential growth**. Later, the slope slows down until we reach the **steady state**: for the SIS the disease keeps circulating among the individuals (*endemic state*), while for the SIR it disappears (*absorbing state*).

In order to compute the **epidemic threshold** for the SIR model, the path to follow is the same as before. In particular, we assume that, at the starting point, $r \ll 1$ so that:

$$\frac{di}{dt} = \beta \langle k \rangle si - \mu i \sim \beta \langle k \rangle i - \mu i \rightarrow i(t) \sim i_0 e^{(\beta \langle k \rangle - \mu)t}$$

and the result we find is again:

$$\beta > \beta_c = \frac{\mu}{\langle k \rangle} \quad (1.12)$$

Since we are able to obtain an analytic expression for S and I in this SIR model, we want to study what is the behavior for large times ($t = \infty$). One obtains that:

$$\frac{ds}{dr} = \frac{-\beta \langle k \rangle s}{\mu}$$

Assuming moreover that $r_0 = 0$ and integrating the above expression wrt r , we obtain:

$$s(t) = s_0 e^{-r(t) \frac{\beta \langle k \rangle}{\mu}}$$

As already said, we cannot find an analytical solution, but we can study the **behavior for large times** by making some approximations. At $t = \infty$, it holds that $i(\infty) = 0$,

thus $s(\infty) = 1 - r(\infty)$ because of the closed population assumption:

$$1 - r(\infty) - s_0 e^{-r(\infty)} \frac{\overbrace{\beta \langle k \rangle}^{R_0}}{\mu} = 0$$

This is a transcendental equation that cannot be solved analytically, but still gives important hints on the behavior of the disease.

One may note that $R_0 = \beta \langle k \rangle / \mu$, and this should make us understand why it is R_0 that drives the exponential growth of the disease, being it proportional to $\beta \langle k \rangle$. Moreover, the initial fraction of susceptible people (s_0) plays a role in shaping the final fraction of recovered. In particular, if $s_0 \ll 1$, the disease cannot spread. This is how **herd immunity** can be obtained.

1.3 Extensions of the SIR model

We want now to modify the SIR to take into account some more features we want to implement our model with.

1.3.1 SIR with Demography

So far we have assumed that the population was totally closed, and so densities always sum up to 1. This is actually unrealistic, so our next step will be to **drop** the **closed population** assumptions: we will now introduce births and deaths (timescale of infection is no more much faster than demographics!). This reasoning is justified from what we observe in real world: considering the demography, we note as every year there are new children that are infected by diseases such as Measles and Chickenpox. Anyway, we do not expect that they will die out over weeks, but still it tells us that newborns increase the populations to the susceptible compartment.

The simplest assumption we can make is: similar to the infectious period, individuals can have a **lifespan**, denoted as $1/\alpha$ years $^{-1}$. Note as in this approximations lifespan is much greater than the infectious period, so *deaths are not due to the disease*. In this way we assume that α is the death rate, common to all classes. Moreover, α is also the crude birth rate, and in addition we assume that births occur only for susceptible individuals and therefore increase its density.

In order to keep the population constant, we need to assume:

$$\frac{ds}{dt} + \frac{di}{dt} + \frac{dr}{dt} = 0 \quad (1.13)$$

Our equations become then:

$$\begin{aligned} \frac{ds}{dt} &= \alpha - \beta si - \alpha s \\ \frac{di}{dt} &= \beta si - \mu i - \alpha i \\ \frac{dr}{dt} &= \mu i - \alpha r \end{aligned} \quad (1.14)$$

where the **infectious period** is:

$$\tau = \frac{1}{\alpha + \mu} \quad (1.15)$$

on average, individuals spend less time infected because some of them may die while infected. However, it is a small change compared to before, since lifespan is much greater than the infectious period.

Also, \mathbf{R}_0 is reduced due to mortality:

$$R_0 = \frac{\beta}{\alpha + \mu} \quad (1.16)$$

We want now to study the **equilibrium points** of the dynamic for this model. Assuming:

$$\frac{ds}{dt} = \frac{di}{dt} = \frac{dr}{dt} = 0$$

we want to find the **equilibrium values** s^* , i^* and r^* . It holds that, at equilibrium:

$$\frac{di}{dt} = 0 = \beta si - \mu i - \alpha i \rightarrow \beta s^* i^* - (\mu + \alpha) i^* = 0$$

and, collecting i^* , we obtain the following equation:

$$i^* [\beta s^* - (\mu + \alpha)] = 0 \quad (1.17)$$

which is not differential anymore.

There are two different solutions for this equation: the one for which $i^* = 0$ (**disease free state**) and the one for $s^* = \frac{\alpha+\mu}{\beta} = \frac{1}{R_0}$, which is the **endemic state**. Here, the most important result is that the **SIR model with demography** can actually **show an endemic state**.

Replacing $s^* = \frac{1}{R_0}$ in $\frac{ds}{dt} = \alpha - \beta si - \alpha w$, we obtain:

$$i^* = \frac{\alpha}{\mu} \left(1 - \frac{1}{R_0} \right) = \frac{\alpha}{\beta} (R_0 - 1)$$

Finally, the three **equilibrium values** (s^*, i^*, r^*) for the fraction of infected, susceptible and recovered in the endemic state are:

$$(s^*, i^*, r^*) = \left(\frac{1}{R_0}, \frac{\alpha}{\beta} (R_0 - 1), 1 - \frac{1}{R_0} - \frac{\alpha}{\beta} (R_0 - 1) \right) \quad (1.18)$$

Keep in mind that this solution exists only if $R_0 > 1$. Moreover, via linear stability analysis, it can be demonstrated that this equilibrium is stable and is reached through damped oscillations.

1.3.2 SIRS Model

We now introduce another model, in which we take into account that during the years the **immune system may lose the ability to recognize a known pathogen**. This immunity could have been acquired via either a vaccine, or having recovered from that disease itself. Moreover, there could be the possibility that viruses mutate, as it occurs with the seasonal influenza, and so antibodies are not able to recognize it any more. Hence, let us build a model in which after an individual is recovered, can become again susceptible after a certain period of time.

The SIRS Model allows to interpolate between SIR ($w = 0$) and SIS ($w \rightarrow \infty$), where w is the **waning immunity rate**, namely the rate at which we lose our ability to defend ourselves from a certain pathogen. We can end up again into either an absorbing, with no more disease, or endemic state, where it keeps on circulating. The transitions for this model are:



Lecture 5.
Thursday 15th
October, 2020.
Compiled:
Thursday 12th
November, 2020.

In particular, the differential equations that describe the model are:

$$\begin{aligned}\frac{ds}{dt} &= \alpha + wr - \beta si - \alpha s \\ \frac{di}{dt} &= \beta si - \mu i - \alpha i \\ \frac{dr}{dt} &= \mu i - wr - \alpha r\end{aligned}\tag{1.20}$$

In this case, the **endemic state** can be found by setting the derivatives equal to zero.

One may note that the transition $R \rightarrow S$ does not affect the I , so it holds that for the **infectious period**:

$$\tau = \frac{1}{\alpha + \mu}\tag{1.21}$$

while the **R_0** factor is:

$$R_0 = \frac{\beta}{\alpha + \mu}\tag{1.22}$$

In addition, the equilibrium values s^* , i^* and r^* can be easily obtained using the same arguments as of the SIR model with demography.

1.3.3 SEIR Model

In reality people do not become instantaneously infectious, but there is a **latent period** which is the time between infection and becoming infectious. Indeed, the pathogen replication takes time, i.e. viral load is too low to be able to transmit the infection. This argument leads us to introduce the **Susceptible, Exposed, Infected, Recovered** model, where the class **E** takes into account that a person has already contracted the disease, hence is not susceptible anymore, but is not able to spread it yet.

Moreover, this period can be extremely heterogeneous depending on the disease: it can take from few hours to years, such as the case for *HIV* or, even longer, *TBC*. In the latter, latent periods might appear to be even longer than an individual's lifespan, with the result that he may have contracted the disease, but the death occurs for other causes before the onset of any symptom.

It is important to remind that the **latent period is not the same of the incubation period** (see Fig. 1.12). An individual can be infectious before symptoms. For instance, there might be a **pre-syntomatic infection period** as it occurs in the case of COVID-19! This explains, once again, why medical status is different from the infection status.

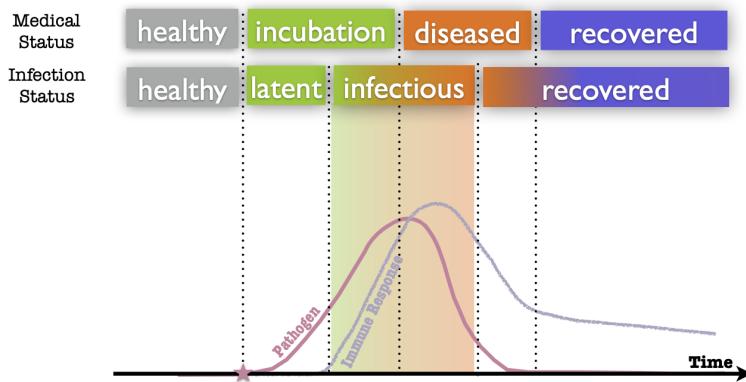


Figure 1.12: Difference between infection status and medical status.

The transition for the **SEIR** model are:

$$\begin{aligned} S + I &\xrightarrow{\beta} I + E \\ E &\xrightarrow{\sigma} I \\ I &\xrightarrow{\mu} R \end{aligned} \tag{1.23}$$

with the equations:

$$\begin{aligned} \frac{ds}{dt} &= \alpha - \beta si - \alpha s \\ \frac{di}{dt} &= \beta si - (\alpha + \sigma)e \\ \frac{dr}{dt} &= \sigma e - (\alpha + \mu)i \end{aligned} \tag{1.24}$$

Hence, the spreading is delayed due to the time spent in E class.

The **endemic state** is:

$$\begin{aligned} s^* &= \frac{(\alpha + \mu)(\alpha + \sigma)}{\beta\sigma} = \frac{1}{R_0} \\ e^* &= \frac{\alpha(\alpha + \mu)}{\beta\sigma}(R_0 - 1) \\ i^* &= \frac{\alpha}{\beta}(R_0 - 1) \end{aligned} \tag{1.25}$$

For very short latent time ($\sigma \rightarrow \infty$) we recover the endemic state of the SIR.

The **R_0** factor is:

$$R_0 = \frac{\beta\sigma}{(\alpha + \mu)(\alpha + \sigma)} \tag{1.26}$$

Since latent time is way shorter than demography one, usually $\frac{\sigma}{\sigma+\alpha} \simeq 1$, hence $R_0 = \frac{\beta}{\alpha+\mu}$ as in the SIR with demography.

One may object that, given that the infectious period and R_0 are similar between SEIR and SIR, adding the Exposed class may seem an unnecessary complication. However, if we look at the time evolution, at the **early stages** there is a huge difference between SEIR and SIR model:

$$\begin{aligned} i_{SEIR}(t) &\approx e^{(\sqrt{4(R_0-1)\sigma\mu+(\sigma+\mu)^2}-(\sigma+\mu))t/2} \approx i_0 e^{(\sqrt{R_0}-1)\mu t} \\ i_{SIR}(t) &\approx i_0 e^{(R_0-1)\mu t} \end{aligned} \tag{1.27}$$

Even if the behavior at the steady state is similar, the temporal evolution of the prevalence of SEIR model is actually slower than the one using SIR. This has surely to be taken into account in policy making, given its important implications.

The SEIR can be the starting point for modeling realistic diseases: i.e. Covid-19 (see Fig. 1.13).

1.4 Summary of compartmental models in well-mixed populations

Let us summarize all the compartmental models in well-mixed populations we have tackled so far:

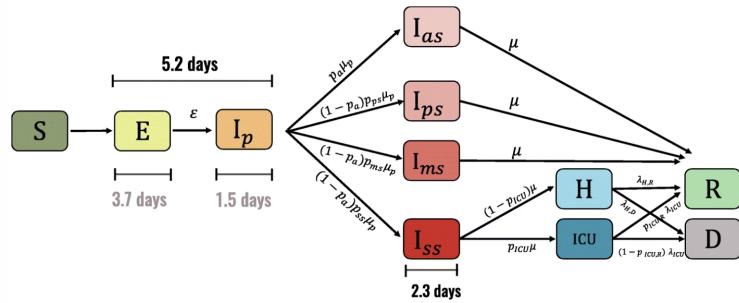


Figure 1.13: Model for Covid-19.

- we solved the **SI model** analytically, and observed that the growth is the one of a sigmoid:

$$i(t) = \frac{i_0 e^{\beta t}}{1 - i_0 + i_0 e^{\beta t}}$$

In the early stages we observe an exponential growth, governed by β , that always saturates at 1;

- in the **SIS model** things starts to change. We have and **endemic** (meta-stable) **state**:

$$i = \frac{\beta - \mu}{\beta}$$

and we reach a sort of **dynamical equilibrium**. We can define an **epidemic threshold**:

$$\beta > \beta_c = \frac{\mu}{\langle k \rangle}$$

- for the **SIR model** equations cannot be solved analytically. However, we observe **no endemic state** and the **epidemic threshold** is once again:

$$\beta > \beta_c = \frac{\mu}{\langle k \rangle}$$

- then, in the **SIRS model** we introduced **waning immunity**. This model interpolates between SIR and SIS model. We do observe **endemic state** and the **infectious period** is:

$$\tau = \frac{1}{\alpha + \mu}$$

and moreover:

$$R_0 = \frac{\beta}{\alpha + \mu}$$

- finally, we discussed **SEIR model** in which we included a **latent period**. We have that:

$$R_0 = \frac{\beta\sigma}{(\alpha + \mu)(\alpha + \sigma)}$$

and the Exposed class has the effect to slow down the spreading.

2

Network Science - Basics

2.1 Main definitions

When we talk about Network Science, as the name would suggest, we study **Networks** that, in math, are also known as graph. A **Graph** $G(V, E)$ is simply an object that is composed by a set of **nodes** (vertices) V and a set of **links** (edges) E :

- **nodes** represent the *entities* $V = [\dots, i, j, k, \dots]$ involved in some relationship. These might be entries, people belonging to a social network and so forth. The **number of nodes** is $N = |V|$;
- **links** represent the relationships between entities $E = [\dots, (i, j), (i, k), \dots]$. The **number of links** is $L = |E|$.

Links can be of different kinds and so networks: the basic distinction is between **undirected** and **directed** links. The former ones can be thought as directed edges, but with arrows pointing in both directions, i.e. to both node of the pair. While the second ones do have a direction according to which sense the relationship represented by the link holds.

Another important distinction is between **unweighted** and **weighted** links. The latter ones can be exploited to take into account the possibility that some nodes can be more connected than the others, therefore **weights** follow. In a certain sense, it describes the "strength" of the link between two nodes.

Another important quantity is the **network density** (connectance), that is the fraction of links present normalized to all the possible pairs:

$$d = \frac{L}{N(N-1)} \quad (2.1)$$

Real networks usually have a very low density, so are **sparse systems** ($L \ll N^2$).

A graph, mathematically, can be represented by the mean of a matrix. It is the so called **adjacency matrix** A of the network, where:

- $a_{ij} = 1$, if a link between nodes i and j exists;
- $a_{ij} = 0$ otherwise.

Many mathematical tools can be used to determine the properties of the system alongside with this matrix, as an example we may want to compute its spectrum in order to obtain the largest eigenvalue. Moreover, one should note that the matrix is symmetrical for undirected and unweighted (also weighted?) graphs, i.e. $a_{ij} = a_{ji}$. However, as we already told, real networks are usually sparse, therefore the adjacency

matrix will be filled for large part by zeros. Hence in order to store graphs in a computer efficiently, it is better to use other tools such as adjacency lists, etc.

Two nodes that share a link are defined "connected", "adjacent", "neighbors". In particular, the **neighborhood** of node i is the set of nodes connected to i . The number of neighbors k_i of each node i is what is called the **degree** of the node i . This is the basic measure that we are going to encounter so many times. Once we have defined the degree, the next step is to define what is the **average degree** over the entire network:

$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i, \quad \text{or} \quad \langle k \rangle = \frac{2L}{N} = d(N-1) \quad (2.2)$$

The next definition is the one of **path**, which is a sequence of links which permits to go from node i to node j following edges. Another relevant quantity is the so called **shortest path** between i and j , it is important since it gives us the idea of how big the network is. In particular, the **distance** l_{ij} represents the length of the shortest path between i and j . There could be multiple shortest paths between i and j . The shortest path of maximum length in the network is defined as **diameter**:

$$l_{max} = \max_{ij} l_{ij}$$

Another measure we may want to introduce is the **average (shortest) path length**:

$$\langle l \rangle = \frac{\sum_{ij} l_{ij}}{N(N-1)}$$

The network is said to be **connected** if every possible couple of nodes is reachable through a path. Otherwise, each connected part is defined as a **connected component**.

Now, let us see some examples of networks, such as "The Oracle of Bacon", or the so called "Erdos Number". The first one is a site that, given the name of an actor, returns the distance between this actor and Kevin Bacon, in unit of costarring movies. This quantity is indeed computed by taking into account the network of actors, linked by common movies in which they starred. The **Erdos Number** instead is the "academical version" for the "Oracle of Bacon": we compute the distance, in terms of collaborations in publications, between a given researcher and the mathematician Paul Erdos through the publications network. The most surprising fact is that, for both examples, the distance is very low! Therefore a question arises: why such short distances in such large networks? In particular, real networks are smaller (i.e. shorter) than one would expect. This is pointed out by the idea of the "Six degrees of separation". It refers to an experiment that was run in the '60s by Stanley Milgram: he gave a postcard to a person on the West Coast, with the instructions that it had to be delivered to a place situated in the East Coast. The main goal was to count how many people would receive that postcard, given the rule that it was allowed to give the postcard to acquaintances of the actual possessor. It was discovered that this postcard actually was delivered to 6 people before reaching the destination. This is what is called the **small world phenomena**. When we study the average path length, for some networks we may find that $\langle l \rangle \sim \log(N)$ or, in some cases even $\langle l \rangle \sim \log(\log(N))$. This is extremely important in the spreading of diseases, since we are able to cover the whole system in few steps.

To summarize what we have seen last lecture: it holds for most real networks that the average path length scales as:

$$\langle l \rangle \approx \ln N$$

Lecture 6.

Friday 16th

October, 2020.

Compiled:

Thursday 12th

November, 2020.

the logarithm of the number of nodes in the network, not just with the number of nodes. Or in some cases as $\langle l \rangle \approx \ln(\ln(N))$. But how is it possible? A paper which explains it is “Collective dynamics of small world networks” by Watts and Strogatz. Their idea is what is called the **Watts and Strogatz model**.

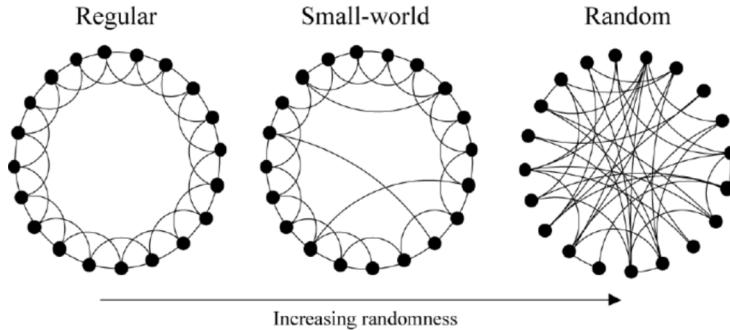


Figure 2.1: Idea of Watts and Strogatz model.

Let us focus on the first regular ring in Fig. 2.1, in which we have that each node is connected with its two nearest neighbours in both sides. The structure as we can see is totally regular. If we want to measure the **longest distance** that we can find in the network:

$$\langle l^{circle} \rangle \sim \frac{N}{4m}$$

But what actually happens if we rewire only a single link? We therefore want to connect it with another random node in the network as in the picture in the middle "small-world" in Fig. 2.1. It can be seen that, by doing a single rewiring, the size of the system reduces in an incredible way. On the other hand, if we extend this argument and choose a probability p for rewiring (i.e. we increase randomness), what happens is that every time we rewire a connection, the average distance is reduced by a factor 2. Repeating this process several times, we observe a **logarithmic scaling**. Finally, the random network we obtain scales as:

$$\langle l \rangle \sim \log N$$

And it is represented by the random circle in Fig. 2.1.

2.2 Degree distribution over networks

Now the question is how degrees are distributed for different type of networks. Let us consider a **small network**, its degree distribution will be really resembling to the plot on the left of Fig. 2.2. However, now we want to understand how this quantity distributes in **real networks**. In order to build a real network, the first assumption that we can make is building the connections *at random*, so with a probability p . Consequently the degree distribution is one of the kind as in the right of Fig. 2.2.

2.2.1 Erdős and Rényi Model: random graphs

Let us consider the Erdős and Rényi model which represents the evolution of a graph where links between nodes are drawn at random, according to a predefined probability p . Before 1959 (the year of the publication of Erdős and Rényi’s paper) people were actually assuming that connections were regular, so no randomness at all. However, since randomness in real world is a deal, thanks to E.R. random connections were taken into account for the first time. In particular, the algorithm for creating such a network is:

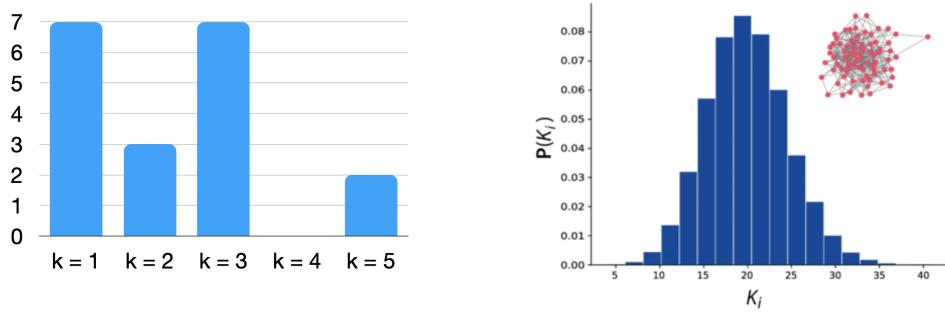


Figure 2.2: **Left:** degree distribution in a small network. **Right:** degree distribution in a network with random connections.

- create an empty graph with N nodes;
- connect each possible couple of nodes with probability p ;
- avoid self-loops and multiple edges.

What are the **properties** of this graph? Let us consider a graph $G(N, p)$, where N are the **number of nodes** and p is the **probability of link**. If links are drawn at random with probability p , the probability p_k that a node has k neighbors is given by a binomial distribution:

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k} \quad (2.3)$$

The **average** and **variance** of such a distribution are:

$$\langle k \rangle = p(N-1), \quad \sigma_k^2 = p(1-p)(N-1) \quad (2.4)$$

As we can see, the average and the variance scales in the same way with the size of the network (i.e. linearly!).

The problem of this distribution is that it is difficult to be dealt with analytically, specially as N increases, indeed:

$$\frac{\sigma_k}{\langle k \rangle} = \sqrt{\frac{1-p}{p(N-1)}} \xrightarrow{N \rightarrow \infty} 0$$

which becomes narrower as N becomes larger, therefore some sort of **approximation** needs to be introduced.

Fortunately, since for sparse networks we have $k \ll N$, the binomial (N, p) distribution can be approximated by a **Poisson distribution** with parameter $\lambda = pN$. Indeed, given that $\langle k \rangle = p(N-1)$, if we have $k \ll N$ then it implies that $p \ll N$. Hence we can write the following:

$$(1-p)^{N-1-k} \approx e^{(N-1-k)\log(1-\langle k \rangle/(N-1))} \xrightarrow{N \rightarrow \infty} e^{-\langle k \rangle}$$

and

$$\binom{N-1}{k} \approx \frac{(N-1)^k}{k!}$$

Obtaining the **Poisson distribution** we were looking for:

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!} \quad (2.5)$$

As before, the average and the variance scale exactly in the same way with the size of the network ($\sim \lambda = Np$). This actually tells us that **all the nodes are more or less the same**. Indeed when we observe a bounded variance, it means that all the nodes more or less have the same degree. In particular, as p increases the graph undergoes a **transition** from disconnected to fully connected one:

- if $Np < 1$, the graph will almost surely have no connected components of size larger than $O(\log(N))$;
- if $Np = 1$, the graph will almost surely have a giant component of size $O(N^{2/3})$;
- if $Np \rightarrow c > 1$, the graph will almost surely have a giant component comprising a large fraction of the nodes;
- if $p < \frac{(1-\varepsilon) \ln N}{N}$, the graph will almost surely contain isolated vertices;
- if $p > \frac{(1-\varepsilon) \ln N}{N}$, the graph will almost surely be connected.

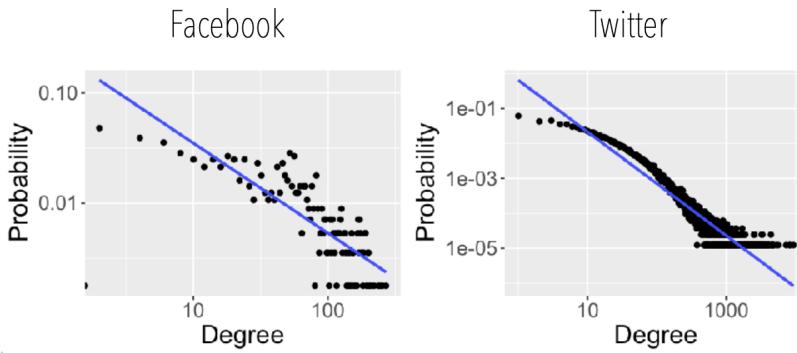


Figure 2.3: Real network of Facebook and Twitter.

2.2.2 Scale-free networks

However, so far we have not discussed how real networks look like, in particular what is their degree distribution. In the last decades we started to have really complex and large networks, whose structure really differs from the structure we usually see for a random network. In Fig. 2.4, as an example, we show two real social networks we know pretty well: Facebook and Twitter. Note that both plots are in log-log scale. Generalizing, we can say that most of the real networks scales in the same way.

We now want to understand how the **degree distribution** looks like. Let us consider Fig. 2.4: black dots follow the Poissonian distribution that we were mentioning before, while the squares follow a power-law $P(k) \sim k^{-\gamma}$, which is **heavy tailed distribution**, in the sense that possibility for large degrees is not null. One should note that the Poissonian distribution is not able to reproduce the heterogeneity we can see in the data, while the power-law is. Hence, in most contexts real networks are **highly heterogeneous** and degrees can span **several orders of magnitude**. In particular, the γ coefficient of the power-law has an important role, since it represents the **slope** of the curve in log-log scale. Since we observe similar structures for different scales, these networks are said to be **scale-free** networks. In most real networks γ has small values, i.e. $\gamma \leq 3$.

Heterogeneity means that almost all nodes have a very low connectivity, way less than a random net. However, the probability of having very large degrees is

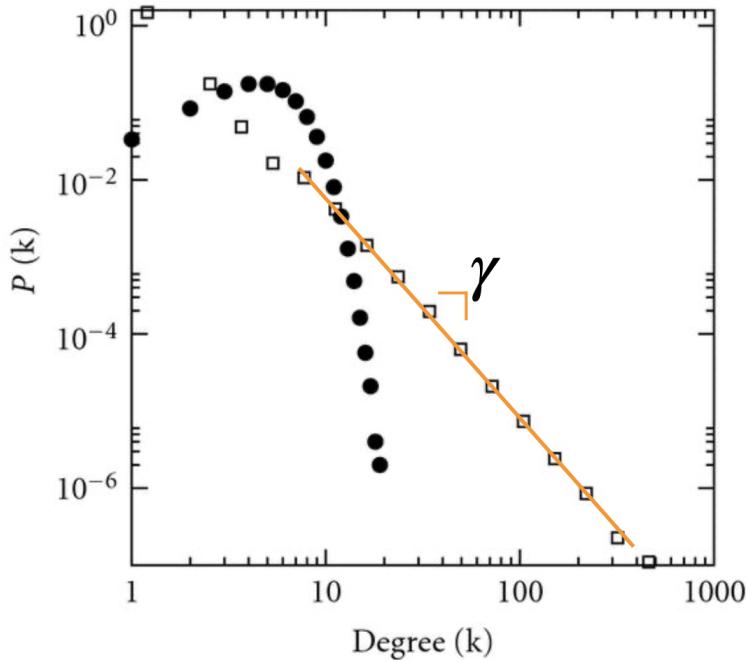


Figure 2.4: Difference between random networks and scale-free networks.

not zero (**hubs**): even for relatively small networks we can observe large hubs. One should take into account that this is something really important for the spreading of diseases: thanks to these large hubs we can see shortcuts for spreading, or the so called **super-spreaders**.

We want now to study the **limiting cases** of these scale-free networks. For instance, we want to see how the **average degree** behaves, or prove that the **largest degree** scales with the size of the network. Let us consider the power-law:

$$P(k) = C_0 k^{-\gamma} \quad \text{with} \quad C_0 = (\gamma - 1) k_{min}^{\gamma-1} \quad (2.6)$$

To understand how k_{max} scales with N , we have to study the case where:

$$\int_{k_{max}}^{\infty} P(k) dk = \frac{1}{N} \rightarrow \left(\frac{k_{min}}{k_{max}} \right)^{\gamma-1} = N$$

Thus, when:

$$k_{max} = k_{min} N^{\frac{1}{\gamma-1}} \quad (2.7)$$

Since in most of networks $\gamma \sim 2 - 3$, so it is easily to understand that k_{max} scales **sub-linearly** with N , but still way faster than random graphs. This is valid for previous plots, such as in Fig. 2.3 as well.

Recalling the definition for the general n^{th} moment of a distribution:

$$\langle k^n \rangle = \int_{k_{min}}^{\infty} k^n P(k) dk = \int_{k_{min}}^{\infty} C_0 k^{n-\gamma} dk \quad (2.8)$$

We note as it converges only if $\gamma - 1 > n$. This gives an hint on how the **average degree** scales as the size of the network: a very important result. If instead we consider the variance $\sigma^2 = \langle k^2 \rangle - \langle k \rangle^2$, we it holds that:

- if $\gamma < 2$, both $\langle k \rangle$ and $\langle k^2 \rangle$ diverge with $N \rightarrow \infty$;
- if $2 < \gamma < 3$, the average degree $\langle k \rangle \rightarrow c$ but $\langle k^2 \rangle \rightarrow \infty$ as $N \rightarrow \infty$, and $\sigma^2 \rightarrow \infty$.

Remembering that most real networks have $\gamma \leq 3$, hence the **variance of the degree also diverges**. The result is that we have extremely **heterogeneous networks** and not homogeneous ones. This is indeed coherent to our observations. It has indeed a very strong **implication**: all the models we have been using before, in which we assumed that **all the people** in the population were **equal**, does **not hold** anymore.

2.2.3 Barabási-Albert Model

So far we have discussed about scale-free networks, but actually we have not created a single one yet. Therefore, an **algorithm** to create such network we can rely on, is the **Barabasi-Albert model**. This topic is discussed in a paper that is the second, chronologically speaking, that gave birth to modern Network Science.

The **idea** behind this paper is extremely simple: once some real networks had been analyzed they assumed that the degree distribution $P(K) \sim k^{-3}$, in order to create a model to reproduce the behaviors observed. Moreover, their model was based on the concept of **growing** for random networks. We start with a small number of nodes, named **clique**, and, at each time-step, a new node enters the network and connects with pre-existing nodes but according to a **preferential attachment**. Therefore, at each step the network grows in size.

The principle on which **preferential attachment** is based on is a very simple concept: *rich gets richer*. That is to say: the more connected a node is, the more likely it is for it to receive new links. The probability for a node i to attract a new link at time t , is proportional to its degree k_i at time t :

$$\Pi(k_i) = \frac{k_i}{\sum_j k_j} \quad (2.9)$$

If we speak about **influencers**, having them a lot of followers, the probability for them to increase their connections is very high. Actually, this idea is not even new, and it is something already known. Indeed this model is just a modification of the *Price model*: if we published a paper and more than someone has found it interesting, it will be more likely for it to receive much more attention in the future.

Specifically for this model, we are drawing links at random, according to some probability that indeed is not uniform. Let us briefly summarize the **main steps** of the algorithm:

- we start with a clique of m_0 nodes;
- at each time step t , we add a new node to the network;
- we create m (i.e. $m = 2$) links between the new node and the existing ones according to the preferential attachment (remember to update the connection probability after each link);
- repeat until the desired size N is reached.

In particular, let us consider Fig. 2.5. We start with a small number of nodes connected via some links. At the *first time step* we add a new node, and then we need to draw connections to the other nodes. Let us assume that every time we add a node, we are adding two links. First, we need to compute the set of probabilities of connecting to each node and, at the first time-step, is equal for all the nodes. Then we pick up one node at random and we draw the link. The following step is to update the set probabilities for each node, according to their degree. We see that the node on the left has got an higher probability of getting new connections since the last node inserted has linked to it. Then, we iterate this procedure by introducing a new

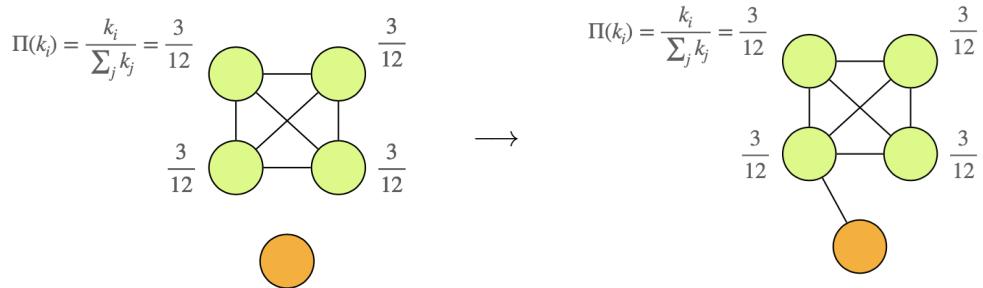


Figure 2.5: Example of Barabási-Albert algorithm.

node and draw connections following the same procedure, until we end up with a total number of N nodes.

This algorithm is indeed able to create networks with some **interesting properties**. Indeed we can approximate the **degree distribution** as:

$$P(k) = \frac{2m(m+1)}{k(k+1)(k+2)} \sim k^{-3}$$

where m is the number of links we are adding at each step. Note that m is a parameter that is related to the minimal degree of the network. However, this approximation is valid for **large** k .

An important result is that $\gamma = 3$ and it is **independent** of m and m_0 . Hence, the **maximum degree** of the network scales as $k_{max} \sim N^{1/2}$. Moreover, it holds that $\langle k \rangle \rightarrow c$, but $\langle k^2 \rangle \rightarrow \infty$ with N , as we have seen before (heterogeneous system). Finally, the **average length** of the network is:

$$\langle l \rangle \sim \frac{\ln(N)}{\ln(\ln(N))}$$

which tells us that the small-world property holds as well.

3

Epidemic Spreading on Networks

Now it is time to drop the assumption of the well-mixed population, and start taking into account **contact networks**. In other words we are considering that **individuals can be connected in different ways** one another. The main idea is that:

- all individuals are **equivalent**;
- we remove the assumption that all individuals have the same number of contacts and we assume that each node **do not interact at random**. This reflects the reality, since we usually have more contacts with some people (friends, family, colleagues...) rather than others. The fact that we may have repeated contacts with someone else has strong effects on the dynamics: we are somehow constraining the way how the disease will spread.

3.1 SIS model in a network

Let us try to build a general model for a general network, without making any assumption on the latter. In order to do that, we start by introducing the equations of SIS model for a generic network.

The first step is to define a **binary variable** for each node i : $\sigma_i(t)$. This variable can only take two values:

- $\sigma_i(t) = 0$, if the individual is **susceptible**;
- $\sigma_i(t) = 1$, if the individual is **infected**.

As one can easily see, this variable describes the state of a generic i -th node at time t . Defining another variable $\rho(i, t)$:

$$\rho(i, t) \equiv \text{Prob}[\sigma_i(t) = 1]$$

which represents the **probability** of that node i is infected at time t . Using this formalism, we can recall the general equation for the SIS in a network:

$$\frac{d}{dt} \rho(i, t) = -\mu \rho(i, t) + \beta \sum_j A_{ij} \text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1] \quad (3.1)$$

The most problematic part is to compute the two nodes infection probability (in green). Since we are in a network, the probability of being infected depends on my neighbours: the (i, j) infection probability depends on the status of all the other neighbors l of j and i and so forth. Therefore we would have to follow the entire

chain of connections, but this would turn out to be a problem: we cannot obtain a closed form for this expression, since it actually depends on the probabilities of all its neighbors. In turn, they would depend on their neighbors probability and so and so forth.

We want to stress one more time that if we want to predict what is going to happen in the system, we would need to consider the entire network and the time evolution for all the nodes. This approach however is **feasible** only for **small graphs** (i.e. 4/5 nodes) and **few compartments**.

This argument reminds us that we may need some sort of an **approximation**: indeed we need to **cut down** this **probability chain**. That is to say that, at some point, we require a closure of our equations, by the mean of approximation: we are not going to take into account the entire structure of the network. At some point we will take the **average**, and after that we will be able to solve the problem. In physics this kind of arguments are called **mean-field approximations**. Since we are not able to solve many body problems, at a certain point we will consider a **random field** which **acts on the entire system** and we will consider its average effects on the system.

Tailoring this procedure to our specific problem, we are substituting in some way the probability $\text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1]$ with some average probability. Obviously, depending on the assumption we are making for this approximation, we will obtain different results.

There are actually many different types of approximations based on different features:

- **Network structure:**
 - **Homogeneous** mean-field (all the nodes are equal);
 - **Heterogeneous** mean-field;
- **Coarsening level:**
 - **Degree-based** mean-field theories (DBMF) in which we assume that all the nodes of the same degree are equal;
 - **Individual-based** mean-field theories (IBMF) in which we assume that all the nodes are different and that we will take individual connections between individuals;
- **Where to cut the chain:**
 - **Individual** level;
 - **Pair** approximations;
 - **Triangles**, etc...;

3.1.1 Homogeneous Networks

Let us start by taking the simplest approximation: we assume **homogeneous network**, **DBMF** and we cut the chain at an **individual** level.

It means that we are considering networks where **nodes degree is bounded**, hence:

- we have that $k_i \simeq \langle k \rangle$;
- we have also that the standard deviation is bounded $\frac{\sigma_k}{\langle k \rangle} = \sqrt{\frac{1-p}{p(N-1)}} \xrightarrow{N \rightarrow \infty} 0$.

All the **nodes can be assumed to be equal**, so their position on the network does not matter anymore. This implies the **spatial homogeneity** it holds that: $\rho(i, t) \equiv \rho(t)$.

In addition, cutting at the individual level means that the two terms of the **joint probability** of one being infected and the other one being susceptible $\text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1] = 0, \sigma_j(t) = 1]$ are **statistically independent**. This implies that the joint probability can be factorized as follows:

$$\text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1] \rightarrow \text{Prob}[\sigma_i(t) = 0] \cdot \text{Prob}[\sigma_j(t) = 1]$$

But now we recall that:

$$\rho(t) = \text{Prob}[\sigma(t) = 1]$$

is the density of infected at time t . Hence, putting everything together, we derive the equation:

$$\frac{d\rho}{dt} = -\mu\rho + \beta \sum_j A_{ij}(1 - \rho)\rho \rightarrow \frac{d\rho}{dt} = -\mu\rho + \beta(1 - \rho)\rho \sum_j A_{ij}$$

Actually, this last term is the degree of the network:

$$\sum_j A_{ij} = k_i \simeq \langle k \rangle \quad (3.2)$$

and by replacing it, we can obtain the same expression that we derived before for SIS model in a well-mixed population:

$$\frac{d\rho}{dt} = \beta \langle k \rangle (1 - \rho)\rho - \mu\rho \quad (3.3)$$

This is a very important result. One should keep in mind that now we are considering all the **nodes statistically independent** and we are back again to exactly the same result of well-mixed population. The only **difference** is that when we were considering well-mixed population, we assumed that the **probabilities** where *exactly statistically independent*. Now, this is just an **approximation**.

Obviously, all the results derived for SIS model in well-mixed populations are still valid, for instance the epidemic threshold.

Remark. Let us recap what we have seen at the end of this lecture. We moved from well-mixed populations to contact networks, so we added more complexity in order to make the model is more realistic. We also derived the equations for SIS dynamics on a generic network and then considered its adjacency matrix. Since for us was impossible to write down a closed equation for this model, given the expression for the infection joint probability that involves two nodes, we were not able to compute exactly the probability for a single node of being infected (ρ_i). It would take into account the probability of three nodes i, j, k at the same time. This is actually unfeasible for all the models and all the possible graphs: it has been done in the literature up to only 4/5 nodes. Hence we end to somehow approximate this probability, in order to cut this infinite chain to a certain value. This is exactly why we introduce mean-field approximation: in this way we take into account the effects of all terms on a specific quantity, not individually, but on average therefore reducing the complexity of our problem. We are switching from a many body problem to a one body problem. The simplest approximation we have seen is the one of homogeneous network in which all the nodes are equal, used on SIS model. According to this argument, for each node there is the same probability of getting infected, so we can approximate the probabilities to be statistically independent. After, we derived all the equations. Their solutions were the same as the ones we had found for well-mixed population. However, in that case the solutions found were *exact*, while now are the result of an approximation.

3.1.2 Heterogeneous Networks

Now, we want to understand what is the effect of **heterogeneity** in the spread of the disease. That is to say that we drop the following assumption $k_i \sim \langle k \rangle$: all **nodes are not equal** any more.

Let us consider now the **heterogeneous mean-field approximation**. Let us use a **DBMF model** and let us cut the chain at an **individual level**. This last assumption means that we consider the probability for a single individual to get the infection. Let us follow the thread of paper “*Epidemic Spreading in Scale-Free Networks*”, written by Pastor-Satorras and Vespignani. It actually provides a **SIS model on scale-free networks**. The main idea behind this paper is the following. Since nodes are not equal anymore, *the probability of getting the infection strongly depends on their position (i.e. degree) in the network*. Authors’ intuition is that **nodes with the same degree behave in the same way**. In order to do that, we need to divide the network in **degree classes**: that is to say we group together all the nodes with the same degree.

In order to write down the equations, we need to consider the number of compartments we have and introduce a density for each of them:

$$s_k = \frac{S_k}{N_k}, \quad \rho_k = \frac{I_k}{N_k}$$

where s_k and ρ_k are the fractions of susceptible and infected nodes of degree k in the network. We have that N_k represents the number of nodes with degree k . As before, we introduced the fractions of susceptible and infected individuals (s_k, ρ_k) in the system, but in this case depending on each degree k . Obviously, the total fraction of ρ and s in the system is given by the sums:

$$\rho = \sum_k P(k)\rho_k, \quad s = \sum_k P(k)s_k \quad (3.4)$$

The **equation** that describes how the **probability of being infected** changes in time for the nodes that belong to the **same degree class**:

$$\frac{d}{dt}\rho_k(t) = -\mu\rho_k(t) + \beta k (1 - \rho_k(t))\Theta_k(t) \quad (3.5)$$

where we can distinct as usual a "recovery" term and an "infection" term. In particular, the probability of a contact between a susceptible individual that has degree k and an infected one is highlighted in green. This product consists in two terms: the probability of being infected $(1 - \rho_k(t))$ and the probability of having contact with an infected $\Theta_k(t)$.

We want now to dwell deeper and explain better this last term. The probability that a generic node with degree k has an infected neighbor can be expressed as:

$$\Theta_k(t) = \sum_{k'} P(k'|k)\rho_{k'} \quad (3.6)$$

where we sum over all the possible degree classes k' . In this way we expect to obtain the probability of connecting with any one of them, multiplied by the probability for that specific node to be infected. Note, however, that we are making no assumption about the function $P(k'|k)$, which may change according k . In principle, it could be anything, in the sense that it strongly depends on the structure of the network. However, in order to simplify the problem and derive some results, there are cases where we can make some assumptions on the structure of the latter.

For **random networks**, e.g. picking a node at random, the probability to be connected to a node of degree k' given the node degree we start from is k , is the following:

$$P(k'|k) = \frac{k'P(k')}{\sum_k kP(k)} = \frac{k'P(k')}{\langle k \rangle} \quad (3.7)$$

where we simply applied the definition of conditional probability. Note that $P(k')$ is the generic probability of getting a connection at random, times k' , which is the number of connection that we pick up (namely the degree k). Finally we normalize over all possible degrees of the network¹. What we obtain is the probability that a generic node in the network is linked to k' . Note as $P(k'|k)$ does not depend on k .

After replacing this last result in 3.6:

$$\Theta_k(t) = \frac{\sum_{k'} P(k')\rho_{k'}(t)}{\langle k \rangle} = \Theta(t)$$

Let us take a look closer to the different terms. In the numerator: there is the product between the probability that a link, randomly picked, points to k' , times the probability of being infected, Finally, we then we sum over all the possible degrees. While, expression on denominator is related only to the structure of the network. In addition, one should note that $\Theta_k(t)$ **does not depend on k** anymore. Since we are just picking up at random it should be the same for all the nodes.

The method that we are going to exploit to **solve** the differential equation $\frac{d}{dt}\rho_k(t)$ is similar to the ones previously used in other models. The first assumption is to be in the **steady state**:

$$\frac{d}{dt}\rho_k(t) = 0 \quad \rightarrow \quad \rho_k = \frac{\beta k \Theta}{\mu + \beta k \Theta}$$

The next step is then to substitute the expression for ρ_k , obtained thanks to Θ :

$$\Theta_k(t) = \frac{\sum_{k'} k'P(k')\rho_{k'}(t)}{\langle k \rangle} = \Theta(t) \quad \rightarrow \quad \Theta = \frac{1}{\langle k \rangle} \sum_k \frac{k^2 P(k) \beta \Theta}{\mu + \beta k \Theta}$$

This is the **self consistent equation** for Θ .

However, in order to solve this last equation, we need some workaround. First of all one should note, as what happens in statistical mechanics, this expression has different solutions depending on the value of Θ :

- the **trivial solution** $\Theta = 0$, that of course is not in our interest;
- the **non trivial solution**. We can rewrite the self consistent equation as follows:

$$\Theta = \frac{1}{\langle k \rangle} \sum_k \frac{k^2 P(k) \beta \Theta}{\mu + \beta k \Theta} = f(\Theta)$$

Hence, the solutions are the values for which it holds $\Theta \equiv f(\Theta)$. These, geometrically, are the intersections between the line Θ and the function $f(\Theta)$ and have to be found graphically (or using computational algorithms).

Since Θ is a probability, it holds that $0 < \Theta \leq 1$. This means that, it is required for a non trivial solution to exist, the slope of $f(\Theta)$ must be greater than 1. Mathematically, it means that:

$$\frac{d}{d\Theta} \left[\frac{1}{\langle k \rangle} \sum_k \frac{k^2 P(k) \beta \Theta}{\mu + \beta k \Theta} \right]_{\Theta=0} \geq 1$$

¹One should keep in mind that $\sum_k kP(k) = \sum_{k'} k'P(k')$.

that leads to the following condition:

$$\frac{\beta}{\mu \langle k \rangle} \sum_k k^2 P(k) \geq 1 \quad \rightarrow \quad \frac{\beta \langle k^2 \rangle}{\mu \langle k \rangle} \geq 1 \quad (3.8)$$

which is the **condition** for the **existence** of an **endemic state**. Since the network has become more complex, also the structure for the condition of the endemic state acquires in complexity. Indeed, for the epidemic threshold:

$$\frac{\beta \langle k^2 \rangle}{\mu \langle k \rangle} = 1 \quad \rightarrow \quad \beta_c = \frac{\mu \langle k \rangle}{\langle k^2 \rangle} \quad (3.9)$$

which is pretty similar to the one previously found, but also includes a term that increases its complexity.

The first check one can make is to verify whether this last result holds also in the case of homogeneous networks. For such networks $\langle k^2 \rangle = \langle k \rangle^2$, therefore:

$$\beta_c = \frac{\mu \langle k \rangle}{\langle k^2 \rangle} = \frac{\mu}{\langle k \rangle}$$

which is exactly the expression we previously found.

Recalling what we were discussing last lectures, in **scale-free networks** with $2 < \gamma \leq 3$, we have $\langle k \rangle \rightarrow c$ and $\langle k^2 \rangle \rightarrow \infty$ as $N \rightarrow \infty$. As the network becomes larger also its variance increases, that is:

$$\beta_c = \frac{\mu \langle k \rangle}{\langle k^2 \rangle} \rightarrow 0$$

hence the **epidemic threshold vanishes** for $N \rightarrow \infty$. This is a quite important result because, if our **network is big enough, every disease will spread, no matter its infectivity** (see Fig. 3.1). The converse is still valid: if we have disease with a very low infection rate in a small part of the network, it will not disappear if the network is large enough²! That is to say we **always** find ourselves in an **endemic state**, while the threshold becomes very small. These results are actually valid for the most real epidemic models, given the networks are large enough.

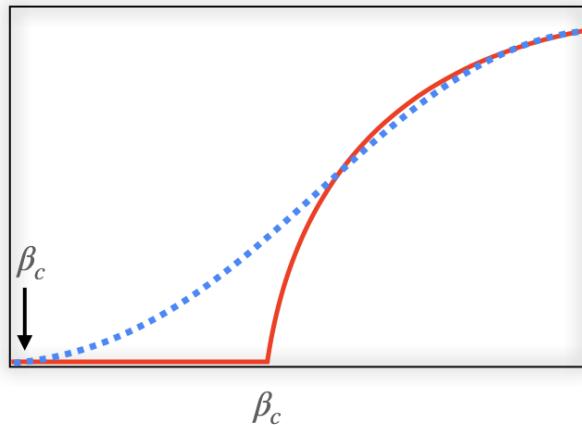


Figure 3.1: In scale-free networks (and many heavy-tailed distributions) the epidemic threshold vanishes in the thermodynamic limit.

²Physically, we refer to this as taking the thermodynamic limit.

Obviously, **real networks** are not infinite: therefore we need some **finite-size corrections**. For example, we may want to derive an expression for epidemic threshold when the size of the system does not diverge.

Let us consider the degree distribution for **scale-free networks**: since the degree cannot go to infinity, it is convenient to introduce an **exponential cut-off** at some point. For instance, let us consider the air transportation network: we see that until a certain point a certain trend is followed, but then the slope of the curve starts to change and resembles to an exponential. This implies that we cannot have an infinite number of connections: the line starts out as a power law and then ends up introducing some sort of exponential cut-off. The behavior is similar to the one in Fig. 3.2.

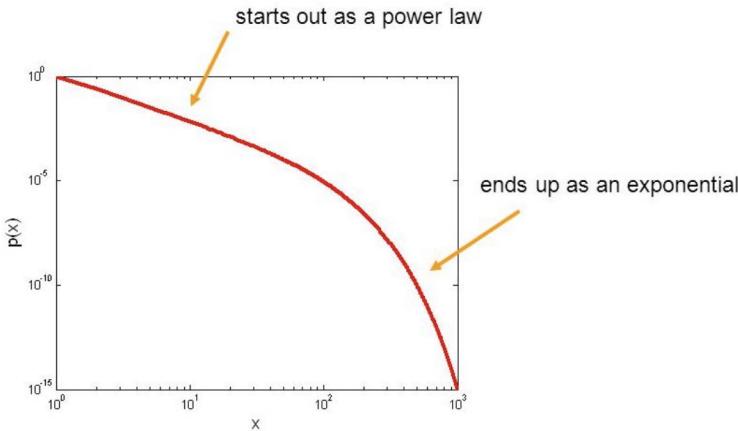


Figure 3.2: Power-law with an exponential cut-off.

We introduce our considerations into our model by adding an exponential term:

$$P(K) \sim k^{-\gamma} e^{-k/k_c} \quad (3.10)$$

where k_c is a **characteristic degree**. At some point, the term we just added will become the dominant term and what happens is that, for large k_c and $2 < \gamma < 3$, the epidemic threshold can be approximated as:

$$\beta_c \simeq \left(\frac{\mu k_c}{k_{min}} \right)^{\gamma-3} \quad (3.11)$$

we are not going to prove the computations. However, in the lab, we will compare the epidemic thresholds for a random and for a scale-free networks in order to see how they differ. This was the last consideration about the study of the SIS model in a network.

3.2 SIR model in a network

3.2.1 Degree-based mean-field theories (DBMF)

The same as before equations can be derived for the SIR model under the assumption of **heterogeneous mean-field**. The main difference is that we need **one more equation** to take into account also the compartment related to **recovered** individuals. Their densities are $\rho_k^S(t)$, $\rho_k^I(t)$ and $\rho_k^R(t)$, and it holds that $\rho_\infty^R = \lim_{t \rightarrow \infty} \sum_k P(k) \rho_k^R(t)$. Equations take the form:

$$\begin{aligned} \frac{d}{dt} \rho_k^I(t) &= -\mu \rho_k^I(t) + \beta k \rho_k^S(t) \Gamma_k(t) \\ \frac{d}{dt} \rho_k^R(t) &= \mu \rho_k^I(t) \end{aligned} \quad (3.12)$$

with $\rho_k^S(t) = 1 - \rho_k^I(t) - \rho_k^R(t)$ and where:

$$\Gamma_k(t) = \sum_{k'} \frac{k' - 1}{k'} P(k'|k) \rho_{k'} \quad (3.13)$$

is the probability of a contact with an infected node, and plays exactly the same role of Θ before. Actually it represents the link from which the infection arrived to that node, however we will not show how to derive this expression beside one small consideration: the $\frac{k' - 1}{k}$ term that is the main difference from the SIS model. It is present due to the fact that we cannot infect a node that has already transmitted us the disease: either because it has already recovered or because it is still infected. In this way we are taking into account that the disease is coming "from one side", therefore for us is forbidden to spread the infection towards that specific direction: recovered (or already infected) individuals cannot be infected twice.

The **epidemic threshold for random networks** results:

$$\beta_c = \frac{\mu \langle k \rangle}{\langle k^2 \rangle - \langle k \rangle} \quad (3.14)$$

and the important thing to notice is that $\beta_c^{SIS} \neq \beta_c^{SIR}$. This is the first time so far that the **epidemic thresholds** for these two models **differ**!

3.2.2 Individual-based mean-field theories (IBMF)

Up to now we were assuming that all the nodes with the same degree were equal. Now, since we are going to study the **individual based mean-field** theories, we will not consider a specific instance of the network, but an average over all the possible networks we can obtain given that **degree distribution**. That is to say, that under the **Heterogenous Mean-Field framework** we are solving the epidemics problem for an **ensemble of networks** whose common feature is the degree distribution $P(k)$ ³.

In the degree based approach we previously assumed that all the nodes with the same degree to be equal. We were therefore analyzing not a specific instance of networks, but its *average*. This is actually what in physics we refer as **annealed networks**. On the opposite, we call **quenched networks** when we consider a *particular realization* of one network. The idea is really simple: instead of considering the average, we consider a particular instance network. This is the main difference between a degree based (i.e. annealed networks) or an individual based approach (i.e. quenched networks).

Let us write down the equations for the **quenched mean-field**. We are going to introduce a **discrete time** framework in order to make equation simpler. However, nothing prevents us to use differential equations, where time is a continuos variable.

Let us consider $\rho_i(t)$ as the probability of a node of being infected at time t . The total fraction of infected individuals is given by $\rho(t) = \sum_i \rho_i(t)$.

At the following time-step, the probability of being infected at time $t + 1$ is:

$$\rho_i(t + 1) = \boxed{\rho_i(t)(1 - \mu)} + \boxed{(1 - \rho_i(t))q_i(t)} \quad (3.15)$$

which is the sum of the probability of being infected and not get cured (green term) and the probability of being susceptible multiplied by the probability of contracting the disease (yellow term).

We now need an expression for $q_i(t)$, that is the **probability** for node i to **be infected**

³the so called "ensemble" of networks!

by, at least, one neighbour. The basic idea for doing this is:

$$q_i(t) = 1 - \prod_{j=1}^N [1 - \beta A_{ij} \rho_j(t)] \quad (3.16)$$

Let us consider Fig. 3.3, in green we have susceptible nodes, which include node i itself, and in red its infected neighbours. The probability of getting infected, at least, by a generic node j is:

$$\beta A_{ij} \rho_j(t) \quad (3.17)$$

Its complementary to 1 is the probability of *NOT* get the infection by node j .

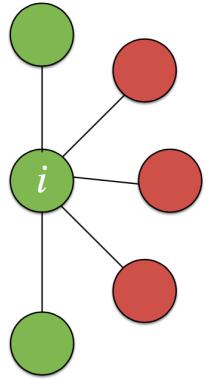
$$[1 - \beta A_{ij} \rho_j(t)] \quad (3.18)$$

Repeating this argument for all neighbors that are actually infected, we can obtain the probability of *NOT* contracting the disease from *ANY* neighbor, namely:

$$\prod_{j=1}^N [1 - \beta A_{ij} \rho_j(t)] \quad (3.19)$$

Again, we previously introduced $q_i(t)$ as the probability of getting infected by at least one neighbor. Hence, the probability of getting infected the complementary to one probability of not getting infected by any neighbor:

$$q_i(t) = 1 - \prod_{j=1}^N [1 - \beta A_{ij} \rho_j(t)] \quad (3.20)$$



Note as the system of ($\rho_i(t+1)$) equations can be solved numerically by iteration. This results to be precise for the entire epidemic diagram, and faster than numerical simulations: there is no need of averages and reproduces individual nodes probabilities. Indeed, in this framework we will obtain two equations for each of the nodes: we have 2^N equations, where N is the size of the system.

Remark. One should have noted that this last approach differs from the degree based mean field theories by the fact that now we are including adjacency matrix A_{ij} , while before we took only the average.

We can also **solve analytically** the system at the **steady state** in order to estimate the **epidemic threshold**. Assuming that we find ourselves in the steady state:

$$\lim_{t \rightarrow \infty} \rho_i(t) = \rho_i^* \quad \rightarrow \quad \rho_i(t+1) = \rho_i(t) = \rho_i^*$$

it follows that:

$$\mu \rho_i^* = (1 - \rho_i^*) q_i^* \quad \rightarrow \quad q_i^* = 1 - \prod_{j=1}^N [1 - \beta A_{ij} \rho_j^*] \quad (3.21)$$

Now, if we think about what happens when we are in **proximity of the epidemic threshold** (*epidemic onset*), it happens that ρ_i^* can be assumed to be small for all the nodes $\rho_i^* = \varepsilon_i^* \ll 1$. Therefore, the product in q_i^* can be approximated by a sum:

$$q_i^* = 1 - \prod_{j=1}^N [1 - \beta A_{ij} \varepsilon_j^*] \simeq \beta \sum_{j=1}^N A_{ij} \varepsilon_j^* \quad (3.22)$$

Figure 3.3: In green susceptible nodes, while in red the infected neighbours.

Substituting what we have just found in the lhs of 3.21 we obtain:

$$\mu \varepsilon_i^* = \beta(1 - \varepsilon_i^*) \sum_{j=1}^N A_{ij} \varepsilon_j^* \quad (3.23)$$

that is a linear system where the interaction is given by the adjacency matrix:

$$\mu \varepsilon_i^* = \beta \sum_{j=1}^N A_{ij} \varepsilon_j^* - \beta \varepsilon_i^* \sum_{j=1}^N A_{ij} \varepsilon_j^*$$

Neglecting second order terms, we have that:

$$\frac{\mu}{\beta} \varepsilon_i^* = \sum_{j=1}^N A_{ij} \varepsilon_j^* \quad (3.24)$$

This linear system has solution only if $\frac{\mu}{\beta}$ is an **eigenvalue** of the **adjacency matrix** A_{ij} . Here we should understand why last lecture we stated that the spectrum of the adjacency matrix is something we may be interested in. Hence:

$$\beta = \frac{\mu}{\Lambda_i} \quad (3.25)$$

where Λ_i is a generic eigenvalue of the adjacency matrix A_{ij} . However, since we are interested in the **smallest** possible **value** of β for which there exists solution, we need to take the **largest eigenvalue** of the adjacency matrix A :

$$\beta_c = \frac{\mu}{\Lambda_{max}} \quad (3.26)$$

The last one is the **expression** for the **epidemic threshold**, and it is a **general result** that is valid not only while using this approximation, but for a more general framework in a generic network.

3.2.3 DBMF vs IBMF: Epidemic threshold

One may wonder now what is the relation between the two values for the epidemic thresholds we have found for the different mean-field theories, that is DBMF and IBMF. We have found that:

- for **DBMF**:

$$\beta_c^{DBMF} = \frac{\mu \langle k \rangle}{\langle k^* \rangle}$$

- for **IMBF**:

$$\beta_c^{IBMF} = \frac{\mu}{\Lambda_{max}}$$

For **scale-free networks** $P(k) \sim k^{-\gamma}$ it holds that:

$$\Lambda_{max} \sim \max \left(\sqrt{k_{max}}, \frac{\langle k^* \rangle}{\langle k \rangle} \right) \quad (3.27)$$

And in particular:

$$\beta_c \sim \begin{cases} \mu / \sqrt{k_{max}} & \gamma > 5/2 \\ \mu \langle k \rangle / \langle k^2 \rangle & 2 < \gamma < 5/2 \end{cases} \quad (3.28)$$

We can conclude that **IBMF** is **more accurate** than DBMF. Due to the approximation, indeed, the **DBMF** is **accurate only** in the **proximity of the epidemic**

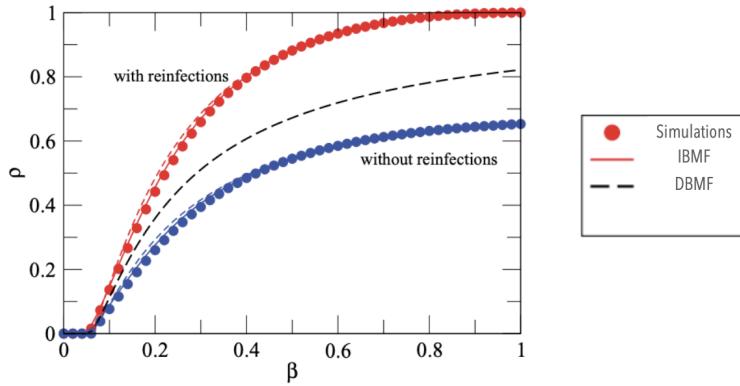


Figure 3.4: The DBMF is accurate only in the proximity of the epidemic threshold while IBMF is accurate for the entire epidemic diagram.

threshold, while IBMF is accurate for the entire epidemic diagram. In Fig. 3.4, this behavior is shown in a clear way. We can see that the quenched mean field follow exactly the simulation, while the DBMF only around the epidemic threshold.

What is the reason behind this result? Since we know the strong connection between the two theories, now we want to **derive DBMF from IBMF**⁴. We start from:

$$\dot{\rho}_i = -\mu\rho_i + (1 - \rho_i)q_i \quad \text{where} \quad q_i = 1 - \prod_{j=1}^N [1 - \beta A_{ij}\rho_j]$$

The Adjacency Matrix is replaced by an **Annealed Adjacency Matrix (AAM)**:

$$\bar{A}_{ij} = \frac{k_j P(k_i|k_j)}{NP(k_i)} \quad (3.29)$$

which for *random networks* becomes:

$$\bar{A}_{ij} = \frac{k_i k_j}{N \langle k \rangle} = \frac{k_i k_j}{2L}$$

where the probability $P(k_i|k_j)$ of picking a random node becomes k_j . The last equation represent the number of trials that we have to create the connection i, j over all the possible connections in the network. If we substitute the explicit form of the adjacency matrix in the expression of q_i :

$$q_i = 1 - \prod_{j=1}^N \left[1 - \beta \frac{k' P(k'|k)}{N_{k'}} \rho_j \right]$$

From individual nodes to degree classes:

$$\dot{\rho}_k = -\mu\rho_k + (1 - \rho_k) \left[1 - \prod_{k'=1}^N \left[1 - \beta \frac{k' P(k'|k)}{N_{k'}} \rho_k \right]^{N_{k'}} \right]$$

Assuming $\beta\rho_k \ll 1$, we can approximate the product with:

$$\dot{\rho}_k = -\mu\rho_k + \beta k (1 - \rho_k) \sum_{k'} P(k'|k) \rho_{k'}$$

Lecture 9.
Thursday 29th
October, 2020.
Compiled:
Thursday 12th
November, 2020.

⁴We repeat that in annealed networks we are not considering a single network but an average of all the possible random network that you can generate from a degree distribution. Instead, in the quenched network, we take a particular one and we want the result for this specific network.

and remembering that $\Theta_k = \sum_{k'} P(k'|k)\rho_{k'}$:

$$\dot{\rho}_k = -\mu\rho_k + \beta k(1 - \rho_k)\Theta_k$$

which is the formula obtained for DBMF. Hence, in the DBMF we are implicitly assuming that $\beta\rho_k \ll 1$. This is the reason why DBMF is accurate only around the epidemic threshold. At the end, we are able to pass from IBMF to DBMF and actually we are explaining the difference in the accuracy between the two models.

3.2.4 IBMF and Pair approximation

Let us make a very brief introduction in what it means cut down the chain to pair approximation. Until now, all the models that we saw were cut at the individual level which means that $\text{Prob}[\sigma_i(t) = 0]$ and $\text{Prob}[\sigma_j(t) = 0]$ are statistically independent:

$$\frac{d}{dt}\rho(i, t) = -\mu\rho(i, t) + \beta \sum_j A_{ij} \text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1]$$

Now, if we cut at the link (i, j) level we have to consider $\text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1]$. This involves three nodes terms and so on:

$$\frac{d}{dt}\rho(i, t) = -\mu\rho(i, t) + \beta \sum_j A_{ij}\rho(j, t) - \beta \sum_j A_{ij}E[X_i(t)X_j(t)]$$

where $E[X_i(t)X_j(t)]$ is the two nodes expectation of being infected. We need an expression for the $\binom{N}{2}$ equations for $E[X_i(t)X_j(t)]$. The idea is:

$$\begin{aligned} \frac{d}{dt}E[X_i(t)X_j(t)] &= -2\mu E[X_i(t)X_j(t)] + \beta \sum_k A_{ik}E[X_j(t)X_k(t)] \\ &\quad + \beta \sum_k A_{jk}E[X_i(t)X_k(t)] - \beta \sum_k (A_{ik} + A_{jk})E[X_i(t)X_j(t)X_k(t)] \end{aligned} \tag{3.30}$$

and the most common possible closures are:

$$E[X_i(t)X_j(t)X_k(t)] = E[X_i(t)X_j(t)]E[X_k(t)]$$

or

$$E[X_i(t)X_j(t)X_k(t)] = \frac{E[X_i(t)X_j(t)]E[X_j(t)X_k(t)]}{E[X_j(t)]}$$

where the second is similar to the first but we are considering the two extremes and then the probability that j is infected.

4

Epidemic spreading on networks: advanced models

In this chapter, we are gonna study non-Markovian epidemic spreading. In the literature, it is not seen, but if you want to implement a realistic model it is very important.

4.1 Markovian Models

All the models we have seen till now assume that β (infection process) and μ (recovery process) are constant rates. This means that movement between compartments takes place at constant rate, or equivalent the probability of jumping from one compartment to another does not depend on the time spent in the compartment. Essentially, we are considering a memoryless process. The jump are memoryless, i.e. we are running a Markov chain.

The fundamental property of a Markov chain is called **Markov property**: *the jump probability at time $t + \Delta t$ does not depend on elapsed time, but only on t* , i.e. we have not to take into account the time that we spent in that compartment. This property is very useful for mathematical treatments.

Jumps made at a constant rate (β and μ) imply that the time spent inside each compartment τ follows a exponential distribution:

$$P(x) = \tau e^{-\tau x} \quad (4.1)$$

with mean $\tau = \frac{1}{\mu}$, i.e. τ is the infectious period (average time spent in I).

What are the implication of this property? The most important one is that an exponentially distributed infectious period implies that the most probable duration of the disease is 0. Indeed, the probability decreases with time. More particularly, it depends on the mean, but in any case the most probable jump (time in which I am making the jump) is at the beginning. It is something that is not realistic, indeed if you got influenza at least you will spend some time infected. And actually, if you are looking how infectius period are distributed in real life, it is something which is quite different. For a disease, you know exactly when it starts but do not you know when it ends. For instance, let us consider the left plot in Fig. 4.1 for 2009 H1N1 Influenza. For this type of disease the plot show the distribtuion of the infectious period, which has as probable value 2 days and an half. The most important things is that it is not 0. On the right, we have also the estimates for Covid-19. One process we use to measure the infectious period is the **serial interval**, i.e. from symptoms to symptom. Or we have also the **generation time**, i.e. from infection to infection. Obviously, these are approximations but can give you some means.

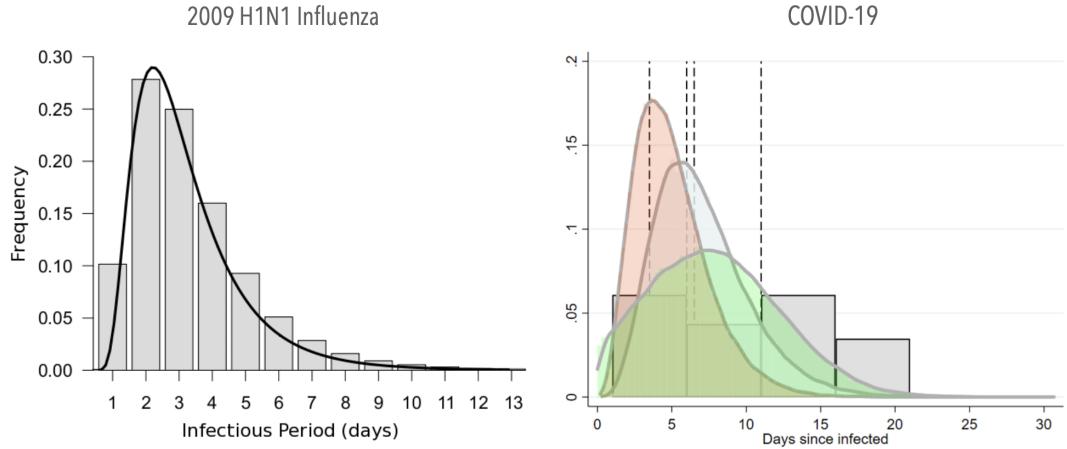


Figure 4.1: Left: Histogram of 1009 H1N1 Influenza. Right: distribution of Covid-19.

The problem is that, all these results demonstrate that these kind of diseases are non-markovian: transition probability depends on the time spent in a compartment. Indeed, patients usually spend some time infected before starting to recover.

4.2 Non-Markovian Epidemic Spreading

How are we gonna model this non markovianity? What are the distribution that better describe what we saw in the data?

Patients usually spend some time infected before starting to recover. This situation is better approximated by a **Gamma distribution**:

$$P(x) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-\frac{x}{\theta}} \quad (4.2)$$

where k is the *shape*, θ the *scale* and $\Gamma(x)$ is the *gamma function*:

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$$

The gamma distribution has mean $k\theta$ and variance $k\theta^2$. This shape start to be somehow what we saw in the data.

What is similar to the gamma distribution is the **Erlang distribution**, where we have the factorial instead of the Gamma function:

$$P(x) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{(k-1)!} \quad (4.3)$$

or the **Weibull distribution**:

$$P(x) = \begin{cases} \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} e^{-(x/\lambda)^k} & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (4.4)$$

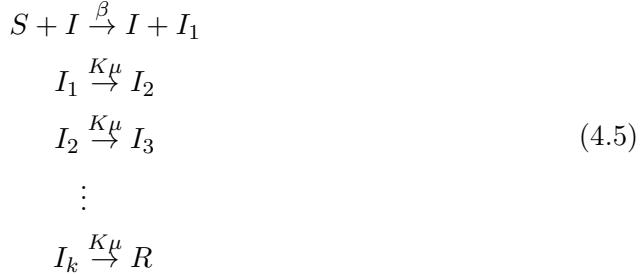
All of these distribution are able to reproduce real histograms.

However, the question still remains: how to include non-markovian elements in classical epidemiological models (with the assumption of markovian...)? We use a trick for the inductive period: *the sum of exponential random variables obeys a Gamma distribution*. How we are gonna incorporate that in our model? Instead of having just one transition at a constant rate (exponential distribution), what we need to have to have a gamma distribution?

The trick is the fact that instead of having just one transition, we are gonna include more and more transitions. Instead of having just one single infectious state, individuals move from one compartment to the other, such that they spend at least some time infectious before starting the recovery. We obtain again a Markovian model.

4.2.1 SIR Model with Multiple Infectious Stages

To repeat, the solution is using **multiple infectious compartment**, i.e. SIR in well mixed populations becomes $S I_1 I_2 \dots I_k R$. For instance, we are imposing that these transitions are sequential:



Hence, if I want to get recovered I need to spent some times infectious, but the model is still markovian! More precisely, the equations are:

$$\begin{aligned} \frac{ds}{dt} &= -\beta si \\ \frac{di_1}{dt} &= \beta si - K\mu i_1 \\ \frac{di_2}{dt} &= K\mu i_1 - K\mu i_2 \\ &\vdots \\ \frac{dr}{dt} &= K\mu i_k \end{aligned} \tag{4.6}$$

where the rate of each I transition is $K\mu$ and with $i = \sum_{k=1}^K i_k$. We got that this is the infectious period distribution:

$$P(\tau) = \frac{(\mu K)^K}{\Gamma(K)} \tau^{K-1} e^{-\mu K \tau} \tag{4.7}$$

where the mean is still $1/\mu$, but the shape is totally different. which is the gamma function. We have two special cases:

- if $K = 1$, we obtain an exponential distribution;
- if $K \rightarrow \infty$ fixed, we obtain a delta distribution.

Other quantities are:

$$R_0 = \frac{\beta}{\mu}$$

and the **final epidemic size**:

$$r_\infty = 1 - e^{-R_0 r_\infty}$$

Moreover, we have an early growth:

$$i(t) \simeq i_0 e^{\lambda t}$$

instead of $i(t) \simeq i_0 e^{(\beta - \mu)t}$. Hence, the disease is growing faster and has a shorter duration.

With λ as the solution of:

$$R_0 = \frac{\lambda}{\mu \left(1 - \left(\frac{\lambda}{K\mu} + 1 \right)^{-K} \right)}$$

The $SI_1I_2 \dots I_kR$ model has several limitaions:

- it is defined only for well-mixed populations;
- focus only on the infectious period distribution. We have that infections are still Markovian and this model only reproduces as Gamma distribution.

4.2.2 Generalized SIS Model

Now, we are going to present something which is more general where we can include non-markovian both in recovery and infections. Is it possible to write down a general model on networks? The answer is yes, it is a bit more complicated and we still needs some kind of approximation at some point. In particular, we need a mean-field approximation.

We have to change our point of view. We are gonna use a slightly different approach, i.e. instead of probabilities we are gonna talk about events. The idea is that we are gonna modelling in this case the infections and recoveries with two random numbers which we extract from distribution and are as general as possible.

The ingredients are:

- a random number $R_i(t)$: recovery time of node i when infected;
- a random number $M_{ij}(t)$: infection times at which node i tries to infect node j .

In order:

- node i get the infection at time t ;
- we extract the random number $R_i(t)$ which represents the time in which node i is gonna recovery (or, the time for which it stay infected);
- then, we extract the random number $M_{ij}(t)$ which represents the number of trials that i try to infect node j while infected;
- we generate a sequence of times

$$T_{ij}^{(1)} \leq \dots \leq T_{ij}^{(M_{ij}(t))} \leq R_i(t)$$

in which node i try to infect node j . For instance, $T_{ij}^{(1)}$ is the first time that node i try to infect node j then we have the second time and so on;

- I am gonna repeat the last step for all my neighbours.

Hence, the transmissibility of the disease is seen as how many trials I am gonna make to infect. One important thing is that $R_i(t)$ and $M_{ij}(t)$ can be drawn from any distribution and not only from the exponential one. How do we extract the T_{ij} is not important at this point, because we are only gonna focus on the distribution of $R_i(t)$ and $M_{ij}(t)$.

Now, let us make some assumptions to make the model more reasonable and then treat it analytically. We assume that:

- $R_i(t)$ and $M_{ij}(t)$ do not depend on time, i.e. $R_i(t) \equiv R_i$ and $M_{ij}(t) \equiv M_{ij}$;
- $R_i(t)$ and $M_{ij}(t)$ do not depend on i and j , i.e. same distribution for all the nodes $R_i \equiv R$ and $M_{ij} \equiv M$.

hence, we are assuming that these numbers should not depend on time and for instance they are typical for the disease. It is valid both for the recovery and for the infections. However, if we consider restrictions as lockdown and so on these number should change, but for the let us consider the simplest model without such restrictions. Indeed, with these assumptions we are reducing the complexity.

We call:

- $E[R]$, the expected value of R ;
- $E[M]$, the expected value of M ;
- v_i , the probability that node i is infected in the steady state.

Now, let us build the model:

1. let us suppose that we are in the steady state of the system (all the transient are passed). In a large time interval $[0, S]$ the number of times node j has been infected is proportional to S (it is linear). Since the length of each infected period is $E[R]$, the **number of infected periods** experienced by a node j (number of times a node has been infected in the interval $[0, S]$) can be written as:

$$\frac{v_j S}{E[R]}$$

2. during each infected period, node j will try to infect i an average $E[M]$ number of times. So, the **total number of infection attempts** from node j to i in a large period of time are:

$$\frac{v_j S E[M]}{E[R]}$$

the number of times j has been infectious multiplied by the number of infection attempts per each time;

3. then, we make the **mean-field assumption**: the conjunct probability that j is infected and i is susceptible is:

$$\text{Prob}[\sigma_i(t) = 0, \sigma_j(t) = 1] \rightarrow \underbrace{\text{Prob}[\sigma_i(t) = 0]}_{(1-v_i)} \underbrace{\text{Prob}[\sigma_j(t) = 1]}_{v_j}$$

4. we can write the **total number of successful infection attempts** for j to i as:

$$S \frac{E[M]}{E[R]} v_j (1 - v_i)$$

that is the total number of attempts multiplied by the probability that i was not infected.

Summing over all the neighbors of i we have the the **total number of successful infections** i will receive during interval $[0, S]$ is:

$$S \sum_{j=1}^N a_{ij} \frac{E[M]}{E[R]} v_j (1 - v_i)$$

5. in the *steady state* the number of successful infections ($S \sum_{j=1}^N a_{ij} \frac{E[M]}{E[R]} v_j (1-v_i)$) is asymptotically equal to the number of infected periods experienced by i : $v_i S / E[R]$. Thus we have:

$$\mathcal{S} \sum_{j=1}^N a_{ij} \frac{E[M]}{E[R]} v_j (1-v_i) = v_i \frac{\mathcal{S}}{E[R]} \quad \rightarrow \quad v_i = E[M] (1-v_i) \sum_{j=1}^N a_{ij} v_j \quad (4.8)$$

hence the probability of i being infected depends by the sum over all its neighbours, times the term $E[M]$ which is the average infection attempts that it is gonna experience during the evolution.

Do the last expression sound familiar? This is exactly the same expression of the IBMF but with a generic infection term $E[M]$:

$$\mu \varepsilon_i^* = \beta (1 - \varepsilon_i^*) \sum_{j=1}^N a_{ij} \varepsilon_j^*$$

The implications are:

- the definition of M implicitly includes the recovery term: $T_{ij}^{(1)} \leq \dots \leq T_{ij}^{(M_{ij}(t))} \leq R_i(t)$;
- exponential case: $E[M]$ is the expected number of infection events in a Poisson process with intensity β within an exponential recovery time with expectation $1/\mu$. Thus: $E[M] = \beta/\mu$.
- epidemic threshold (lower bound):

$$m_c = E[M_c] = \frac{1}{\Lambda_{max}}$$

Hence, this is the form in which we can generate a generic model from.

Today we will take into account something which is quite common. Usually diseases does not spread in isolation, but interact with other diseases.

Most of the ifnection diseases that we have seen write now, there are different variants circulating for the same disease. For instance, for the season influenza we have different virus which are similar, we have a family of virus. These virus interact between each other.

It is extremely important that a disease can be with a several variants.

Let us consider in which we have two different diseases which cooperate. One disease boost the spread of the other one. The map represent the spread of tuberclosis in 1990. Tuberclosi is a disease with a very strong latent period, the disease can be in a latent state for many many years. We have the same map in 2005. We note that the disease exploit especially in the south. The number almost doubles. The reason is that HIV reach africa. HIV is a disease which compromise your immune system, so your immune system is less efficient: the probability of getting any other disease is higher than the normal.

When they got HIV the immune system go down and tuberclosi activate. The numbers represent the fraction of pations which get HIV and tuberclosis. SO HIV is helping a lot for the spread of tuberclosi. This was the first example of interaction of two diseases.

The other case we can have is **competition** as seasonal influenza. We have the distribution of the different influenza virus which circulates for the season. We have a cake plot in which we see that a certain number of people got this type of influenza

Lecture 11.
Thursday 5th
November, 2020.
Compiled:
Thursday 12th
November, 2020.

(light blue). We can see that the influenza curve is a mixture of several different types.

Each year seasonal influenza is caused by different virus.

Our immune system has a sort of memory. So if we got one disease we can get the same the next year.

So actually the susceptible population is just a reduced fraction. This is why every year we do not have a huge pandemic of influenza.

There is another case **evolution** when we have virus which mutate a lot. For instance HIV is an extremely volatile virus so it can mutate a lot. We have different types and subtypes. We can see how these subtypes are distributed on different areas.

We are gonna focus only on the simple setting in which we consider competition and cooperation between diseases. How is it possible to model these kind of interactions? The simplest solution is to couple different dynamics.

For instance, in the simplest state we can couple SIS and SIR dynamics. We have twice the number of states (we take into account all the possible combinations). Each disease its gonna have its own parameter. We need at least some parameter which encode this kind of interaction between diseases.

Let us see one of this model. Let us consider a classical heterogeneous mean-field in which we have two different diseases. We have also two different networks in which diseases spread. For instance, one disease diffuse in an oral way, while the other one by blood. So the network can be different from one disease to the other. So, each network has its own degree distribution and we have to use the joint distribution.

Since we are doing a degree based mean field, in this case we have four different equations for all couple of different disease. So, the state of the system is represented by these four variables.

Hence, we have to take into account (to represent a general model) this interaction effect. We introduce the **modified susceptibility**. For instance, having one disease make the individual more susceptible than the others.

If we want that this framework as general as possible, we have to consider that the suscpebility can be modiefied and also the probability of getting the diseas. If I am being infected of one disease makes me more/less infected (depending if the parameter is positive or not).

The third interaction that you can have is on the recovery. If I have another disease probably my recovery time will be large. If η is smaller than 1, my recovery period will be larger (cooperative effects of the diseases), instead if it is larger than 1 we have the competitive.

These are all the possible interactions we can have in this simple model. These are all the possibilities that we can have. The normal ifneciion and the recovery process. Then, we have an increase in suscpetibility, in infectivity. You can play with these parameter as much as you want. The model cover all the possibilities but only a subset of them cover a biological sense.

We talked about the flows of variation. How people changes inside the compartments. (FORMULATION)

The structure is absolutely the same, also the way of solving it. I am not showing all the results, some passages are extremely large. We assume that we are in the steady state, then we write down a self-consistent equation for σ_1 and σ_2 and we solve exactly in the same way by finding the intersection. At the end we got this expression for the epidemic treshold:

$$\beta$$

this is very very complicated but has the same form of before.

The epidemic treshold of the first disease depends on the prevalnce of the other. I am assuming that one disease is already there and I am putting the second disease

there and I am seein what are the effect for the epidemic treshold.

Let us see what happens in different cases. For instance, the epidemic diagram is in 3D for coopereting disease with $\lambda > 0$ and $\eta < 1$.

Weh I put only one disease we obtain the classical results (in 2D). Since I am not putting the second disease here we have no disease. Whwn we consider two disease, untile the second disease is below the treshold, I am not seeing any different (the disease is spreading as before), but then the epidemic treshold is gonna decrease. I gonna have a larger presence of the disease.

Things works exactly in the same way when I have competition between the two disease for $\lambda < 0$ and $\eta > 1$. Both of the diseases cannot spread and I have no disease in that area.

This was the case for the heterogeneous mean field degree based. We can write also the quenched mean.field equation (individual based formulation).

We consider just one network. We consider only the effect of modified susceptibility. Also in this case we can write the equations for the quenched mean field. These are the equations. The structure is exactly the same. Each term is one contribution to this probability $[\rho^{IS}]_i^{t+1}$.

We insert the function f just to consider the fact that If I am getting two disease, we pick only one disease at random.

We can solve this problem numerically, by iteration as we were doing fot the single case scenario. Let us see what happens when we have cooperation between diseases. We have the classical case with $\lambda = 1$. When the probability of getting the disease is $\lambda = 2$ we see.. by increasing we see that we see the same sort of exploiting behavior and the curve becomes more deeper. My infectiviti is exacly the same, but the fraction of nodes is two times, three times... we see large jumps.

Let us see the full cross-imminuity: if you have one disease you cannot get the other. If we plot the total fraction of infected this is not chaning at all (is exactly the same of before). Instead of plotting the sum of them, if we plot the difference we see that after one point the difference is exactly the same value of the same meaning that only one disease can survive. Since the diseases are symmetrycal the one which survive are only by chance. After some infectivity what happens is that you can see the difference going to zero. Both of them survive and we have some sort of equal range. One half of the population for one disease and one half for the other. For large beta the two diseases coexisti.....

Actually, we can also think in terms of the dynamics of the system. The system have two stable points: when there is on disease and the other one is zero. AFter a treshold the stable point is exacyl in the middle. Both diseases are coexisting with exactly the same incidence.

When you have a slight difference, the point in the middle its gonna move on this line.

5

Spreading in social systems

We want to show how we can apply the kind of epidemic models described in the previous chapters to other scenario, especially in social systems. Indeed, spreading of informations on social systems share many similarity with epidemic spreading (i.e. “viral information”). There is a huge literature of models adapted to this kind of thing and we can adapt the models we have seen until now to include social system aspects.

However, there are also some differences; since the communication aspect is different we have some effects that are not totally included in simple epidemic models. In social contexts things are a bit more complex:

- information transmission is an *intentional* act for both sender and receiver;
- often beneficial for senders and receivers (e.g. **reinforcement**). Indeed, the information can be replicated by different sources (for instance, when I see an information different times in different places);
- influenced by cognitive and psychological factors;
- content of information matters (e.g. **homophily**).

For all these reasons this type of spreading is usually defined as:

- **simplex contagion** (there is no memory, no reinforcement) against the **complex contagion**. Complex contagion may involve multiple exposures and reinforcement, i.e. memory of past interactions vs independent interactions.
- we call the **threshold models** since we see some sort of threshold effects (if half of my friend buy an Iphone I will likely buy it too).

In a normal model, when you are susceptible, your infected neighbours try to infect you with a certain probability. But, in this case, what is happening is that if your neighbours number is lower than a certain threshold, you cannot be infected by them, while if they are equal or above you are gonna change your state.

In particular, threshold models lead to information cascades.

5.1 Complex contagion

We want to have some sort of general model which is able to interpolate between complex and simple contagion.

A generalized contagion model able to reproduce both simple and complex contagion is constituted by:

- we are gonna assume to be in well-mixed population;

- we consider S (susceptible), I (infected) and R (recovered) individuals;
- we want memory of past interactions up to time T (we want to include the fact that we have different exposures in different period of times.);
- then, we change the way in which the information is spread. From each successful interaction with an infected j , a susceptible i gets a “dose” of infection $d_i(t)$ (If I see an information more time the probability of spreading it will be higher);
- if the accumulated dose $D_i(t) = \sum_{t'=t-T+1}^t d_i(t')$ exceeds a threshold d_i^* , then the susceptible becomes infected.

These are modifications we can make to the SIR dynamics. Regarding the **infection** process:

- at each time step t :
 - each individual i contacts a random individual j ;
 - if $i = S$ and $j = I$, with probability p individual i gets a “dose” of infection $d_i(t)$ from a dose size distribution $f(d)$;
 - otherwise (if the contact is not successful) $d_i(t) = 0$.
- each individual keeps a “cumulative” dose of the T previous time steps: $D_i(t) = \sum_{t'=t-T+1}^t d_i(t')$.
- if $D_i(t)$ is larger than the individual threshold d_i^* , individual i gets infected.

For the **recovery** process is more or less like the classical dynamics:

- if $D_i(t)$ gets below d_i^* , i recovers with probability r ;
- it also possible to add an $R \rightarrow S$ transition with probability r' to simulate an SIRS model (reinfection dynamics). With $r = 1$ and $r' = 1$ it recovers an SIS-like dynamics.

From that, let us summarize the main parameters that we have:

- p and r infection and recovery probability (same role as β and μ);
- $d_i(t)$ “dose” per infection distributed following $f(d)$;
- d_i^* threshold, distributed following $g(d^*)$.

Note that $f(d)$ and $g(d^*)$ can be general distributions. By changing them we can reproduce different behaviors (different kind of dynamics). With specific choices of p , $f(d)$ and $g(d^*)$, it is possible to recover both pure epidemic and threshold dynamics.

We can formulate the model mathematically and we are gonna try to solve it. First of all, let us define the probability that an individual with $K < T$ contacts gets infected as:

$$P_{inf}(K) = \sum_{k=1}^K \binom{K}{k} p^k (1-p)^{K-k} P_k \quad (5.1)$$

where $p^k (1-p)^{K-k}$ is the probability that the contact is successful (this is the Bernoulli distribution of having K trials and k success). This is multiplied by all the possibilities $\binom{K}{k}$ and then for each contact we multiply by the probability P_k .

In particular P_k is the average fraction of infected after receiving k doses in T time steps:

$$P_k = \int_0^\infty dd^* g(d^*) P\left(\sum_{i=1}^k d_i \geq d^*\right) \quad (5.2)$$

which obviously depends on our threshold. Indeed, $P\left(\sum_{i=1}^k d_i \geq d^*\right)$ is the probability that k doses exceeds d^* .

The model can be solved numerically for any distribution of $f(d)$ and $g(d^*)$. For some specific cases we can recover classical dynamics. Let us consider:

- if $p < 1$ and fixed dose size $f(d) = \delta(d-1)$ and fixed threshold $g(d^*) = \delta(d^*-1)$ we have **epidemic spreading** (independent interactions). In particular, all contacts have the same infection probability and the threshold is $d^* = 1$ (one successful contact). See Fig. 5.1. Hence, if we want to recover an SIR dynamics, we fix that each dose is exactly the same for everyone (picked around one) and also the threshold should be one;

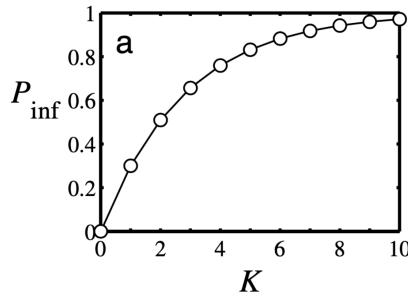


Figure 5.1: Generalized Complex Contagion model: epidemic spreading.

- if $p = 1$ and fixed dose size $f(d) = \delta(d-1)$ and fixed threshold $g(d^*) = \delta(d^*-5)$ we have **deterministic threshold model**. In particular, the threshold is fixed at $d^* = 5$, so we need at least 5 encounters to get infected. Hence, the dose size is exactly the same of before (since the distribution is still peaked around one), but I need more than one contact (for instance 5 friends show to me the information). See Fig. 5.2;

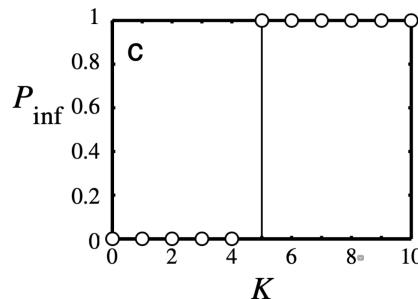


Figure 5.2: Generalized Complex Contagion model: deterministic threshold model.

- if $p = 1$ and fixed dose size $f(d)$ is distributed log-normally and fixed threshold $g(d^*) = \delta(d^*-5)$ we have **stochastic threshold model**. In particular, the threshold is still fixed at $d^* = 5$, but the “dose” for each contact varies. Hence in this case we are considering that the contacts are not equal (for instance I trust one friend more than another) and we get that the threshold is of the same, but the dose size is different. See Fig. 5.3.

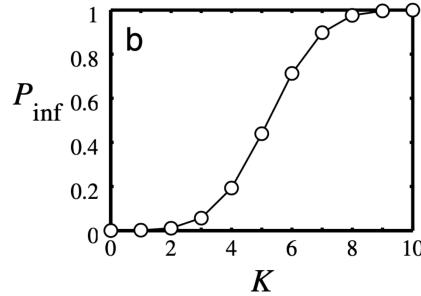


Figure 5.3: Generalized Complex Contagion model: stochastic threshold model.

5.2 Applications to Online Social Networks

The question is: how hashtags (memes) spread in online social networks? Let us consider the analysis of real data on Twitter. In twitter we have different types of communitiess. We are gonna study of retweets (RT) and mentions (@) diffuse hashtags in different communities and in particular we want to see the effect of communities (i.e. how reinforcmenet and homophily will have a role on spreading)

To quantify the fraction of information that flows inside a community and outside, we are gonna use the following weights:

- $\langle w_{\circ} \rangle_c$ is the average weight (number of tweets) per link inside the community;
- $\langle w_{\sim} \rangle_c$ is the average weight (number of tweets) per link outside the community.

The same for users activity:

- f_{\circ} is the fraction of activity inside the community;
- f_{\sim} is the fraction of activity outside the community.

If information in Twitter spread like a *simple contagion* there should be no differences in spreading inside and outside a community (e.g. *no reinforcement*) Instead, if we see that the average weights inside the community are major than outside actually you need some sort of reinforcement.

In Fig. 5.4 we can see the results showing the average weight inside a community and outside. They are pretty similar, but what happens is that it seems that the average are a little bit higher for spreading inside a community (homophily and reinforcement have a role).

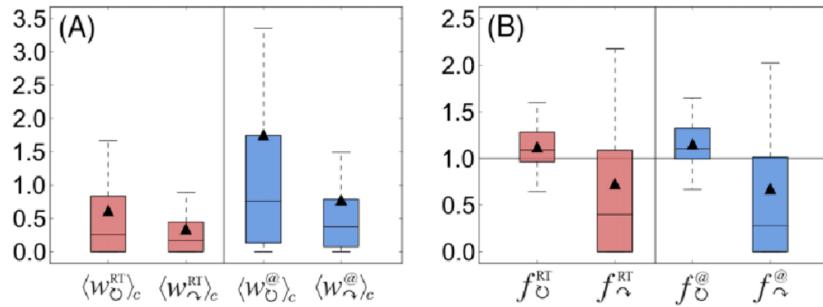


Figure 5.4: Community spreading is favored (effects of homophily and reinforcement).

To clarify that point we introduce more measures at the simple hastag (h) level: we measure the hashtag inside a community and outside (i.e. for each hastag we measure the popularity that a tweet got). In particular, for each hastag:

- we measure the **usage dominance** $r(h)$: proportion of tweets produced inside the “main” community of h out of the total number of tweets containing h , $T(h)$. We expect that this measure is low for viral spreading and high for complex contagion;
- we measure the **usage entropy** $H(h)$: how h is distributed across communities. It is high for viral spreading and low for complex contagion;
- we measure the **average exposure** $N(h)$: average number of exposures needed to adopt hashtag h . It is low for viral spreading and high for complex contagion.

We have some reference models (4 models $M_{1,\dots,4}$) to represent different baseline behavior (see Fig. 5.5 for more details):

- the simplest one is when for a given hastag h , the model M_1 randomlu samples the same number of tweets or users as in the real data (i.e. data are extracted at random). We got some sort of **average behavior** for all the hastags (no community, no network structure);
- the second model is just an epidemic model. In particular, M_2 takes the *network structure* into account while neglecting social reinforcement and homophily. Each hastag starts with a random users and at each step with a certain probability spread to the other user. This is a reference model for **simple contagion**;
- we can have more complex model which takes into account the *network structure, reinforcement* and *homophily*. This is a reference model for **complex contagion**.

Table 1 Baseline models for information diffusion			
Community effects			Simulation implementation
Network	Reinforcement	Homophily	
M_1			For a given hashtag h , M_1 randomly samples the same number of tweets or users as in the real data.
M_2	✓		M_2 takes the network structure into account while neglecting social reinforcement and homophily. M_2 starts with a random seed user. At each step, with probability p , an infected node is randomly selected and one of its neighbors adopts the meme, or with probability $1 - p$, the process restarts from a new seed user ($p = 0.85$).
M_3	✓	✓	The cascade in M_3 is generated similarly to M_2 but at each step the user with the maximum number of infected neighbors adopts the meme.
M_4	✓		In M_4 , the simple cascading process is simulated in the same way as in M_2 but subject to the constraint that at each step, only neighbors in the same community have a chance to adopt the meme.

Average behavior
Simple contagion
Complex contagion

Figure 5.5: Reference models for information spreading in online social networks.

Let us consider Fig. 5.6 where we can see $r(h)$, $H(h)$ and $N(h)$ as function of the number of tweets T and the number of users U . The black lines represent the real data, the dashed lines represent model M_1 (average behavior), the red square M_2 (simple contagion) while the blue and green models M_3 and M_4 (complex contagion).

There is a real distinction between popular (in grey) and not popular hastags. We see two main behaviors:

- popular hashtags (large T and U) spread like epidemics (viral);
- less popular ones follow a complex contagion.

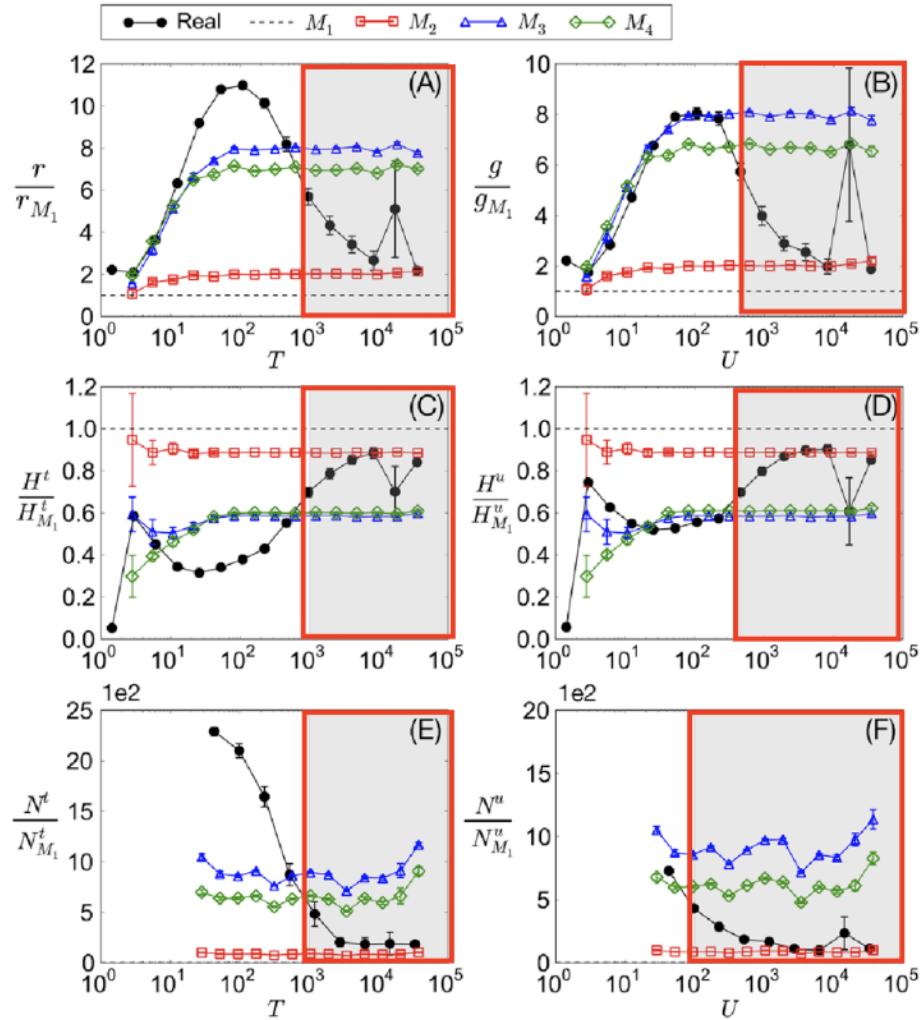


Figure 5.6: Results of information spreading in online social networks.

Part II

Poletto's Lectures

Bibliography