# A Spatial Markov Chain Cellular Automata Model for the Spread of Viruses

Jenny Lu[1] and Fred Vermolen[2,1(✉)]

[1] Delft Institute of Applied Mathematics, Delft University of Technology,
Mekelweg 4, 2628 CD Delft, The Netherlands
`J.Lu@student.tudelft.nl`

[2] Computational Mathematics Group, Department of Mathematics and Statistics,
University of Hasselt, Agoralaan Building D, 3590 Diepenbeek, Belgium
`Fred.Vermolen@uhasselt.be`

**Abstract.** In this paper a Spatial Markov Chain Cellular Automata model for the spread of viruses is proposed. The model is based on a graph with connected nodes, where the nodes represent individuals and the connections between the nodes denote the relations between humans. In this way, a graph is connected where the probability of infectious spread from person to person is determined by the intensity of interpersonal contact. Infectious transfer is determined by chance. The model is extended to incorporate various lockdown scenarios. Simulations with different lockdowns are provided. In addition, under logistic regression, the probability of death as a function of age and gender is estimated, as well as the duration of the disease given that an individual dies from it. The estimations have been done based on actual data of RIVM (from the Netherlands).

## 1 Introduction

At the time of writing, the COVID-19 crisis has been affecting the global human population worldwide for almost two years. The impact of the virus has been enormous, almost five million people have died [1], and many countries have been in various degrees of lockdowns and economies have been hit hard. The lockdown policies have caused countries to close their borders, ban travel, people to live in isolation and companies to suffer great losses. Sadly, despite the vaccinations, the pandemic is far from over. There are new mutated versions of the virus, like the beta, gamma, delta and omikron versions. In addition, recovered people can lose their immunity against the virus in a couple of months, hence making reinfection a problem as well.

The disease, COVID-19, is renamed as the strain Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) by the World Health Organisation [2]. It is often characterised by flu-like symptoms, which in some cases lead to excessive fever or even to lung inflammations. One of the serious problems regarding

this disease is the high infection rate from person to person. In addition, the disease affects every individual. Even in young people, COVID-19 can cause strokes, seizures and Guillain-Barre syndrome—a condition that causes temporary paralysis. COVID-19 may also increase the risk of developing Parkinson's disease and Alzheimer's disease according to the Mayo Clinic [3]. The three organs that are impacted most by the disease are the heart, lungs and brain. The long term effects of the disease are yet still unknown.

In order to predict the dynamics of the spread, death rate and recovery rate of the coronavirus, many different strategies are used. A very common model is the so-called SIR model, see [4] for the original paper. More modern elaborations on the SIR principle have been presented in [5,6] and [7]. In particular, the model in [7] bears some similarities with the approach that is presented in the current paper. This model simulates a homogeneous population that is exposed to a virus. It contains a susceptible, infected, recovered and dead fraction of the population. Many more advanced models are variations on this strategy. One attempts to include spatial spread by the incorporation of diffusion terms, which are justified by random (unpredictable) migration and interaction of individuals. Other extensions are based on the incorporation of networks, which allows so-called jump processes so that airborne communication can be taken into account. The models described in [8] distinguish several challenges for network modelling in epidemics. A review on mathematical modelling of epidemics has been written in [9]. This review considers the different practices and limitations of modelling global spread of diseases. Duan et al. [10] wrote a review about epidemic modelling where models of different nature are discussed. First the S(E)IR models based on ordinary differential equations are introduced, and this is followed by network models that are based on stochastic principles. It is reported that stochastic (network) models yield very realistic results [10]. Stochastic network models for epidemics have been presented in [11], where an exact final size distribution is constructed using recursive formulas. Further, the impact of vaccination is quantified in [11]. A Bayesian inference for stochastic epidemic models has been presented in [12]. Both [12] and [11] favour the use of stochastic network models because of the huge flexibility from temporal and random effects that the models are able to handle. The current model elaborates on the influence of the topology of the network on the evolution of the epidemic. The results that are presented in the current paper should be classified as preliminary in the sense that the results are based on simulation with hypothetical input parameters. However, later in the paper we do present an estimation for the recovery rate parameter based on actual data of the National Institute for Public Health and the Environment or short RIVM of the Netherlands. The infection parameter is hard to estimate based on the data, due to the many factors that influence the possibility of infection and because of the implementation of measures that change over time.

The model presented in this paper is based on cellular automata, in which the nodes of the grid represent individuals who are connected to each other by means of a graph. Each individual is assigned a state at every time instance, these states

are: susceptible, infected, recovered or dead. The stochastic nature of the model can be seen by the probability of infection as well as death. Contact between individuals does not always lead to infection, and hence here a stochastic process is considered. The probability that an individual infects another is determined by the intensity of the contacts that the individuals have. Lockdown policies have been implemented in the model by adjusting a pre-specified parameter. Next to being infected, recovery is incorporated and once an individual has recovered, then it is assumed that the individual is immune to the disease. We are aware of possible reinfection, and we study this topic, but for the current manuscript, reinfection is omitted. Since COVID-19 can be a lethal disease in some cases, death has been incorporated as well. The model has been extended for modelling lockdown policies that certain governments have adopted. One of the advantages of the current approach is the small number of input parameters needed. A further innovation is the uncertainty quantification and the statistical assessment of the results.

## 2   The Mathematical Model

In this section the mathematical model will be derived and explained. First a basic model is given, which is later extended to a more realistic model.

To begin, consider a graph with nodes and vertices. The nodes represent individuals that can be in one of four different states: susceptible, infected, recovered, dead. If a person is susceptible, then this individual can only be infected. Once the individual is infected, then, the person can either recover or die. If (s)he recovers then this person is assumed to be resistant. If a person is susceptible, dead or resistant, then (s)he will not spread the virus to other people (although this assumption may be subject to discussion because a non-infected could spread the virus via the hands or other objects, however, this effect is neglected in the current modelling). The interpersonal relations are represented by connecting line segments in the graph. The connection is subject to an intensity, which represents the frequency that two individuals physically interact. This intensity and connection can be interpreted in a generic sense regarding relations and geographical distances. This connection determines the probability that, if one of the two individuals is infected, the disease is transferred from one another. Furthermore, infected individuals may recover or die.

Mathematically, this can be written as follows. The population consists of $n$ individuals, which at every time instance is denoted by a vector of length $n$, where each entry in this vector contains the state of individual. This vector is denoted by $\mathbf{v}$, where the value of $v_i$ contains the integer states: $v_i \in \{1, 2, 3, 4\}$, where $v_i = 1$, $v_i = 2$, $v_i = 3$ and $v_i = 4$, respectively, correspond to the susceptible, infected, recovered and dead states. All individuals are connected to other individuals by vertices between nodes (or individuals). The connection between person $i$ and $j$ is denoted by $a_{ij}(t)$, where $a_{ij}(t) = 0$ represents the case that individuals $i$ and $j$ have no physical contact. The entries $a_{ij}(t)$ are assembled into the contact intensity matrix $A(t)$. Note that the entries are dependent on time $t$ as the

intensity of the contact between two individuals changes over time. Large values of $a_{ij}(t)$ represent a high intensity of the contacts, while lower values mean that the people have less contact. The dynamics of the spread of the disease is discussed in the coming subsections.

## 2.1   The Transfer of the Virus from Individual to Individual

Suppose that individual $i$ is infected and that individual $j$ is susceptible. The time between going from the susceptible state to the infected state is assumed to follow an exponential distribution with infection parameter rate $\lambda_{ij}(t)$ in a given time interval denoted by $\tau$. The time interval $\tau$ is assumed to be small. Hence, the probability that person $j$ becomes infected in the small time interval $\tau$, given that person $i$ is infected, is given by:

$$P(v_j(t+\tau)=2|v_j(t)=1, v_i(t)=2) = \int_t^{t+\tau} \lambda_{ij}(s)e^{-\lambda_{ij}(s)(s-t)}ds. \quad (1)$$

Since it is only possible to go from the susceptible state to the infected state, the probability that person $j$ stays susceptible in a small time interval $\tau$, given that he/she was susceptible at time $t$ and person $i$ is infected at time $t$, is given by:

$$P(v_j(t+\tau)=1|v_j(t)=1, v_i(t)=2) = 1 - \int_t^{t+\tau} \lambda_{ij}(s)e^{-\lambda_{ij}(s)(s-t)}ds. \quad (2)$$

Therefore, the probability that this non-infected person $j$ dies or recovers from the disease is zero. Hence:

$$P(v_j(t+\tau) \in \{3,4\}|v_j(t)=1, v_i(t)=2) = 0. \quad (3)$$

The infection rate parameter $\lambda_{ij}(t)$ is assumed to be of the following form:

$$\lambda_{ij}(t) = \lambda_g \, a_{ij}(t), \quad (4)$$

where $\lambda_g$ is a general infection rate parameter that is assumed to be the same for every individual.

Next we consider the set of people an individual is in contact with. Define the set $N_j$ of individuals that is in contact with person $j$ by:

$$N_j(t) = \{k \in \{1, ..., n\} : a_{kj}(t) > 0\}, \quad (5)$$

where the set $N_j(t)$ represents the 'neighbours' of person $j$; these are the of individuals are in contact with person $j$ at time $t$. This set can be reduced to a set where we only consider all the neighbours of individual $j$ that are in the infected state. This subset is denoted by $N_j^I$:

$$N_j^I(t) = \{k \in N_j(t) : v_k(t) = 2\}, \quad (6)$$

where the superscript $I$ denotes the infected individuals.

Next, assume that all the states of the individuals in the 'neighbour' set $N_j$ or $N_j^I$ are independent of each other. Hence to obtain the probability of not being infected, the product of all the probabilities can be used. Therefore, the probability that node $j$ stays susceptible is as follows:

$$P(v_j(t+\tau) = 1 | v_j(t) = 1) = \prod_{k \in N_j^I(t)} (1 - \int_t^{t+\tau} \lambda_{kj}(s) e^{-\lambda_{kj}(s)(s-t)} ds). \qquad (7)$$

As a direct consequence, the probability that node $j$ becomes infected is given by:

$$P(v_j(t+\tau) = 2 | v_j(t) = 1) = 1 - \prod_{k \in N_j^I(t)} (1 - \int_t^{t+\tau} \lambda_{kj}(s) e^{-\lambda_{kj}(s)(s-t)} ds). \qquad (8)$$

During the time interval $[t, t+\tau]$ where $s \in [t, t+\tau]$ we take $\lambda_{kj}$ constant, hence $\lambda_{kj}(s) = \lambda_{kj}(t)$. Then Eq. (8) can be rewritten as:

$$P(v_j(t+\tau) = 2 | v_j(t) = 1) = 1 - \prod_{k \in N_j^I(t)} e^{-\lambda_{kj}(t)\tau} = 1 - e^{-\tau \sum_{k \in N_j^I(t)} \lambda_{kj}(t)}. \qquad (9)$$

Substituting the definition of $\lambda_{kj}(t)$ gives:

$$P(v_j(t+\tau) = 2 | v_j(t) = 1) = 1 - \prod_{k \in N_j^I(t)} e^{-\lambda_{kj}(t)\tau} = 1 - e^{-\tau \sum_{k \in N_j^I(t)} a_{kj}(t)\lambda_g}. \qquad (10)$$

From equation (10), it can be seen that there is an effective transfer probability rate that can defined for each node $v_j$, namely:

$$\lambda_j^{eff} = \lambda_g \sum_{k \in N_j^I(t)} a_{kj}(t). \qquad (11)$$

This is summarised in Theorem 1, of which a similar version was proved in [13]

**Theorem 1:** *Let node $i$ possess neigbours $\mathcal{N}_i^I(t)$ that are infected. Then, assuming the contact intensity matrix not to change during the time interval $(t, t+\tau)$, the effective probability rate in the exponential distribution for node $i$ to become infected is given by*

$$\lambda_i^{eff} = \lambda_g \sum_{j \in \mathcal{N}_i^I(t)} a_{ij}(t).$$

## 2.2   Transition to Recovery and Death

People that are in the infected state await two different scenarios: recovery with being resistant or death. Some people recover very quickly after having had (very) mild or even no symptoms, whereas other people need a long time to recover or pass away. In the current modelling, it is assumed that the recovery time follows an exponential distribution with probability rate parameter $\mu > 0$, that is

$$P(v_i(t + \tau) = 3|v_i(t) = 2) = \int_t^{t+\tau} \mu e^{-\mu(s-t)} ds. \tag{12}$$

It has been assumed that $\mu$ is constant, which is not realistic as $\mu$ is be subject to temporal changes due to improvements of medical therapies against the disease as well as the health conditions of each individual or even seasonal effects. Later this assumption is relaxed. The probability that a person stays infected is then given by:

$$P(v_i(t + \tau) = 2|v_i(t) = 2) = 1 - \int_t^{t+\tau} \mu e^{-\mu(s-t)} ds. \tag{13}$$

The expected recovery time from the moment that the patient was infected is determined by

$$T_r = \frac{1}{\mu}, \tag{14}$$

which follows from the properties of the exponential distribution. In this model it is assumed that if a person has been infected during a time-interval that exceeds a threshold, say $T_{death} = M \times \tau$, where $M > 0$ is some positive integer value, then the person dies with probability one. This is not always the case in reality, but it is a reasonable assumption. The rationale behind this assumption is that a long lasting exposure to the disease potentially damages the patient's vital organs so much that he/she dies.

To develop our intuition behind the relation between the recovery rate and time interval of death $T_d$, we assume that the probability that someone dies from the disease is given by $\alpha$. Hence all patients that have been ill over a period that exceeds $T_d$ are assumed to die. Then if person $i$ was infected at time $t$, then the time interval of death and the probability to die are related by:

$$1 - \int_t^{t+T_d} \mu e^{-\mu(s-t)} ds = \alpha \tag{15}$$

We can rewrite it as:

$$\int_0^{T_d} \mu e^{-\mu(s-t)} ds = 1 - \alpha \implies T_d = -\frac{1}{\mu}\log(\alpha) = -\log(\alpha)T_r,$$

where log is the natural logarithm and $\alpha$ is a very small probability that is at most 2–3 %.

If it is assumed that person $i$ got infected at time $t_{inf} = \min_{t>0}\{v_i(t) = 2\}$ we can write the following mathematically:

$$v_i(t_{inf} + \theta) = \begin{cases} 3, & \text{if } \theta < T_d \\ 4, & \text{if } \theta \geq T_d \end{cases} \tag{16}$$

where $\theta$ is amount of time that person $i$ has been infected.

## 3 Computational Implementation

In all the simulations, a constant population size is assumed. The current preliminary computations involve a simplified square topology, in which each node has at most four connections. It is easy to revise this topology. A uniform transmission probability rate $\lambda_g$ to obtain the probability that the node changes from susceptible to infected during the time step $\tau$. To model transmission, the effective transmission probability rate is computed by the use of the contact intensity matrix. Subsequently for each susceptible node a random number, $\xi$, from the standard uniform distribution (between zero and one) is sampled, that is $\xi \sim U(0, 1)$. If the number is smaller than the probability of transmission from susceptible to infected then the state is changed from susceptible to infected, that is $v_i$ is changed from 1 to 2, that is

$$v_i(t + \tau) = \begin{cases} 2, & \xi < P(v_i(t + \tau) = 2 | v_i(t) = 1), \\ \\ 1, & \xi > P(v_i(t + \tau) = 1 | v_i(t) = 1). \end{cases}$$

Otherwise, it stays in the susceptible state.

The same is done for the transition from the infected state to the resistant or dead state. However, we keep track of the time-interval that a node has remained in the infected state by adding the time step $\tau$ to the time-interval that a node is in the infected state. If the total length of the time interval that a nodal point stays in the infected state exceeds the length $T_d$, then the node is moved to the dead state. As long as this time-interval has not been exceeded, the node either stays in the infected state or it is transferred to the resistant state analogously to the treatment of susceptible nodes, but with a different number $\xi_2 \sim U(0, 1)$ drawn from the uniform distribution.

Since interpersonal contacts are often fluctuating (like going to shops, meeting friends, working, etc.), randomised values for the contact intensity matrix $a_{ij}(t)$ are used, that is, considering person $i$:

$$\text{For } s \in (t, t + \tau) \ : \ a_{ij}(s) \sim U(0, 1), \qquad \text{if } j \in \mathcal{N}_i.$$

In the case of lockdown, the contact intensity matrix is multiplied by a factor, $\beta(t)$, whose value ranges between zero and one. Small values of the factor $\beta(t)$ represent severe lockdown policies. Hence, the contact intensity matrix is redefined by

$$\hat{A}(t) = \beta(t)A(t),$$

and in all expressions given earlier, $A(t)$ and its entries are replaced with $\hat{A}(t)$ and its corresponding entries are given by $\hat{a}_{ij}(t) = \beta(t)a_{ij}(t)$. Note that the lockdown policy depends on $t$, and therefore $\beta = \beta(t)$, where $\beta : \mathbb{R}^+ \longrightarrow [0, 1]$.

In order to compute the fractions of susceptible, dead, infected and resistant people, we introduce the standard Kronecker Delta Function:

$$\delta_{p,k} : \mathbb{N} \times \mathbb{N} \longrightarrow \{0,1\} \ : \qquad \delta_{p,k} = \begin{cases} 1, & \text{if } k = p, \\ 0, & \text{else.} \end{cases}$$

The fraction of individuals in state $p \in \{1,2,3,4\}$ is given by

$$f_p(t) = \frac{1}{n} \sum_{j=1}^{n} \delta_{p,v_j(t)}.$$

This definition reproduces that $\sum_{p \in \{1,2,3,4\}} f_p = 1$. The computations are termi-
nated as soon as the number of infected people equals zero.

## 4 Simulation Results

In the previous sections the mathematical model has been derived and the numerical implementation has been explained. In this section the results of various simulations will be shown and examined.

In all of these simulations, a rectangular arrangement of $100 \times 100$ nodes is taken. Every internal node will only have four neighbours (left, right, up, down), every boundary node will have three neighbours and every corner node will only have two neighbours. Initially all the nodes are in susceptible state, except for a node in the lower left corner, that node is infected (this is a random choice, one could have made any node in the population infected). Without an infected person, the model will never predict the spread of the virus. The reason why it has been chosen in the lower left corner is that this grid can be seen as one of four 'quadrants', so the results can be reflected towards the other three quadrants. It is to be noted that the results of the current simulations are hypothetical as we have only used hypothetical values. Later on in this paper, parameter estimation will be done in order to find the input parameters of the model, given a simulated data set. This enables an estimate of the input parameters based on observed data and therewith the model can be made predictive (Table 1).
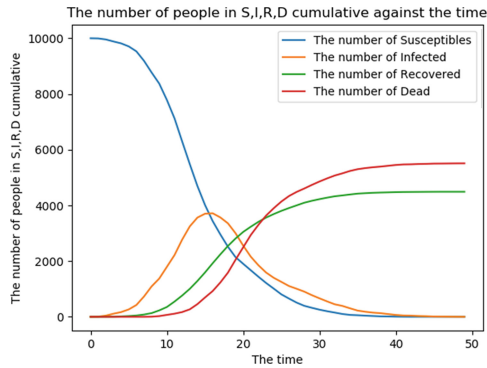
### 4.1 Simulation 1: No Lockdown

We begin by examining the scenario where there is no lockdown, hence in the model this means that $\beta(t) = 1$. It is expected that the virus will be able to spread rapidly in the population, causing many active cases in a short period of time. The virus rapidly spreads within the population until eventually everybody has either remained susceptible, recovered or died. We will refer to the different

**Table 1.** Table of parameters used in the simulations

| Parameter | Value |
|---|---|
| Grid $(n_x \times n_y)$ | $100 \times 100$ |
| $\tau$ | 1 |
| $\lambda_g$ | 0.5 |
| $\mu$ | 0.1 |
| end time | $100\tau$ |
| Time lockdown start | $15\tau$ |
| Time lockdown end | $60\tau$ |
| Time to death $T_d$ | 8 |

groups of people in the population as the *Susceptible, Infected, Recovered* and *Dead* sub-populations. These sub-populations consider the number susceptible, infected, recovered and dead people respectively at a certain time interval in the simulation. Capital letters have been chosen as we refer to this specific sub-population group. To get a better understanding of how the virus spreads, the sub-populations have been plotted against the time in Fig. 1, the Recovered and Dead sub-populations have been plotted cumulative against the time. The Susceptible and Infected sub-populations have not been plotted cumulative, because the number of susceptible people only decreases over time and eventually all the infected people recover or die.

Figure 1 shows that the virus spreads exponentially in the beginning, causing a rapid growth in the number of active cases and a drop in the number of susceptible people. The Infected sub-population graph looks almost like a bell-curve shape. It starts with exponential growth up until time $15\tau$ and then it makes a round turn and has a more parabola shape afterwards, until it flattens out as there are no more infected people left.



**Fig. 1.** Graph of the number of people in Susceptible, Infected, Recovered, Dead sub-populations cumulative against the time with no lockdown.

The reason why there is a turn in the number of infected people is because at time $15\tau$ we do not have 'enough' susceptible people left to obtain a higher peak than before, therefore the number of infections decreases afterwards. The Susceptible sub-population graph has a very steep negative slope, suggesting that the number of susceptible people rapidly decreases. The number of infected people follow a similar trend, but then reversed (so a rapid increase). After the number of infected people have peaked, the decrease in the number of susceptible people also slows down, until there are no more susceptible people remaining. The cumulative graphs of the Recovered and Dead sub-populations in Fig. 1 look like logistic growth. The Recovered sub-population graph starts to rise earlier compared to the Dead sub-population graph, as it takes some time for infected people to recover. In this situation there are more dead people than recovered people, but that is due to the way the parameters in the simulation have been chosen. If a different set of parameters were chosen, the outcome would be different.

**Uncertainties in the Simulations when there is no Lockdown**

The mathematical model that is presented is based on random processes, therefore each simulation will have a different graph of all the four sub-populations. To see the possible 'bandwidth' that a sub-population might have, the simulation has been carried out a hundred times and the various sub-populations have been plotted against the time. These graphs are found in Fig. 2.
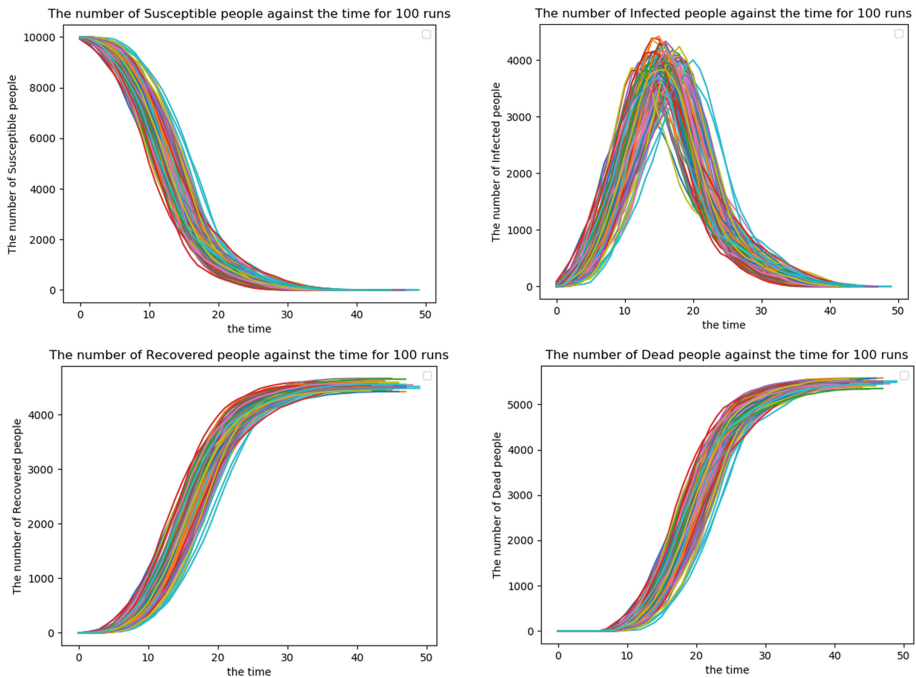


**Fig. 2.** Cumulative graphs of the number of people in recovered and dead sub-populations against the time with no lockdown for 100 simulations. The graphs of the susceptible and infected sub-populations are not cumulative.

If we consider the top two figures in Fig. 2 more closely, then it is noticeable that the peak of the Infected sub-population graphs is around time $15\tau$ if there is no lockdown implemented. This also corresponds to the steepest slope in the Susceptible people sub-population graphs, suggesting it might be interesting to incorporate a lockdown scenario at time $15\tau$. This will be done later in this section.

From Fig. 2 it can be seen that the number of recovered people fluctuates around the 4000 to 4500 (when considering a total population size of 10000 people), which is roughly 40–45% of the total population. The number of dead people ranges between 5000 and 5500, which is roughly 50–55% of the total population. This is of course not the ideal outcome, but it is the result of the parameters chosen in the simulations. Figure 2 shows that all the simulation results are consistent, but possess some variation.

## 4.2   Applying a Temporary Lockdown

In this section some simulations will be shown where the lockdown is not kept constant over time or lifted immediately, but there is a step in between. First a case is considered, where a severe lockdown of $\beta(t) = 0.1$ is implemented, then loosened to $\beta(t) = 0.5$ and finally to $\beta(t) = 1$ (no lockdown). In addition, another case is considered where first a heavy lockdown of $\beta(t) = 0.3$ is implemented and then lifted to a lockdown of $\beta(t) = 0.6$ or $\beta(t) = 0.7$ to see what impact such lockdown has. The parameters that are used in the simulations are found in Table 2.

**Table 2.** Table of parameters used in the simulations

| Parameter | Value |
|---|---|
| Grid $(n_x \times n_y)$ | $100 \times 100$ |
| $\tau$ | 1 |
| $\lambda_g$ | 0.5 |
| $\mu$ | 0.1 |
| End time | $100\tau$ |
| Time lockdown 1 start | $15\tau$ |
| Time lockdown 2 start | $35\tau$ |
| Time lockdown end | $60\tau$ |
| Time to death $T_d$ | $8\tau$ |

**Case 1: Severe Lockdown of $\beta(t) = 0.1$ to Medium Lockdown of $\beta(t) = 0.5$ and then no Lockdown**
We start by examining a simulation where first a severe lockdown of $\beta(t) = 0.1$

has been implemented. This lockdown is later lifted to a medium lockdown of $\beta(t) = 0.5$ at time $35\tau$ and at time $60\tau$ it is entirely lifted. Therefore:

$$\beta(t) = \begin{cases} \tilde{\beta}(t) = 0.1, & \text{for } t \in T_{ld_1} = (15\tau, 35\tau) \\ \tilde{\beta}(t) = 0.5, & \text{for } t \in T_{ld_2} = (35\tau, 60\tau) \\ 1, & \text{else.} \end{cases}$$

This could refer to the case where a country has first implemented a strict protocol where people are only allowed to go outside for one hour a day and all public restaurants/events/bars/shops are closed. Then the rules are lifted to a medium lockdown where shops are open again but only for limited amount of customers and everybody has to wear a face mask, and at time $60\tau$ everything is back to normal again. At time $100\tau$, when the simulation ended, there is still a significant amount of susceptible people left as well as infected people. This indicates that the virus is still spreading amongst the population. In order to understand what consequences this lockdown implementation has brought to the population, consider the graphs of the various sub-populations in the first figure in Fig. 3. From this figure some interesting events as a result of the lift of the lockdown rules can be seen. Until time $15\tau$ everything is like before as the virus is free to spread. At time $15\tau$, the consequence of the severe lockdown of $\beta(t) = 0.1$ can be seen (just like in Case 3 of the previous subsections). Due to the relaxation of the lockdown at time $35\tau$, a second outbreak occurs, which causes a rise in the number of infections. However, the number of susceptible people still decreases at a relatively constant rate, until time $60\tau$. At time $60\tau$ the lockdown is lifted and we notice a third peak in the number of active cases. The second and third peak are nevertheless relatively small compared to the first peak. This is not only a result of the lockdown rules, but also a result of less susceptible people in the population. The most striking observation in the Recovered and Dead graphs is that the third peak is not really visible, suggesting that the relaxation of a medium lockdown to no lockdown does not affect the number of recovered or dead people significantly. However, this might be explained by the population size. If a larger population size would be taken (say 100 million) this would be more significant.

**Case 2: Heavy Lockdown of $\beta(t) = 0.3$ to Mild Lockdown of $\beta(t) = 0.6$ and then no Lockdown**

The next scenario is as follows: a country has first implemented a heavy lockdown of $\beta(t) = 0.3$, then relaxed the rules to a mild lockdown of $\beta(t) = 0.6$ and after that the lockdown is entirely lifted. The results can be seen in the second graph in Fig. 3.

$$\beta(t) = \begin{cases} \tilde{\beta}(t) = 0.3, & \text{for } t \in T_{ld_1} = (15\tau, 35\tau) \\ \tilde{\beta}(t) = 0.6, & \text{for } t \in T_{ld_2} = (35\tau, 60\tau) \\ 1, & \text{else.} \end{cases}$$

The effect of the first lockdown of $\beta(t) = 0.3$ is significantly present as can be seen in the number of active cases. The loosening of the rules at time $35\tau$ causes another rise in the number of infections from time $35\tau$ until time $42\tau$. It is quite

interesting to see that the first peak looks almost triangular while the second peak looks like half an oval (shape wise). This is probably due to the fact that the first peak is caused by no lockdown rules, so we first have an exponential growth of infected people. Then suddenly a heavy lockdown is implemented which causes the amount of infected people to drop drastically. The second lockdown is mild, which causes a rise in the number of active cases, but since there are few susceptible people remaining in the population, it quickly turns and starts decreasing again. The peaks in the number of active cases correspond to the rapid increases in the Recovered and Dead sub-population graphs. In addition, another simulation has been run, but in this case, the second lockdown was lifted till 0.7 instead of 0.6 to see if this makes any difference. The results can be seen the third graph in Fig. 3. The difference is not significant since the same trends as in the second case are observed.
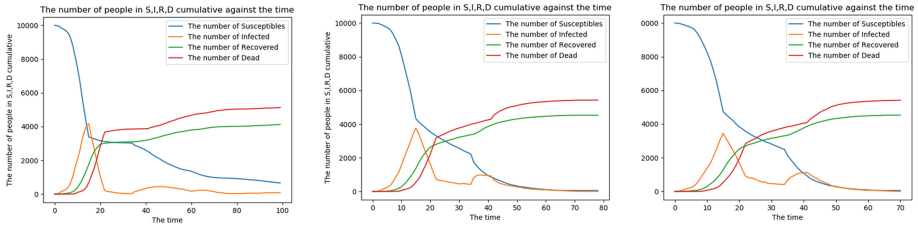


**Fig. 3.** Cumulative graph of the sub-populations against the time. The first graph refers to case 1, where $\beta(t) = 0.1$ to 0.5 to 1.0, the second to case 2, where $\beta(t) = 0.3$ changes to 0.6 to 1.0, the last graph refers to $\beta(t) = 0.3$ to 0.7 to 1.0.

## 5    Estimation of Recovery and Death

In this section an estimate for the recovery rate (hence also the death rate) and an estimate of the duration of COVID-19, given that the patient dies, are provided. The estimates are done by the use of statistical models, such as logistic regression for the probability of recovery as a function of age and gender, linear regression, log-linear regression and log-Poisson regression for the duration of the disease given that the patient dies as a function of age of the patient. Similar approaches have been carried out by Alleman et al. [6] and by Chen [14]. The estimations are based on the data provided by the Dutch Institute for Health and Environment (RIVM). First some data processing is done, followed by the underlying mathematical model and the parameter estimation. The estimation of the infection parameter $\lambda_g$ is not done based on real data due to the fact that it is so highly dependent on the lockdown degree (which was changed almost every two weeks in the Netherlands), the testing rate, the number of contacts, age, vaccination etc.

## 5.1   Data Processing

The data provided by the RIVM consists of a couple of variables, defined as follows:

– the date that the person reported that he/she got infected by COVID at the GGD (Municipal Health Services of the Netherlands).
– the gender of the patient, which in this case is Male, Female or Unreported.
– the age of the patient in years.
– first date that the patient experienced symptoms of COVID.
– the date that the patient died because of COVID.

There are a few things to be noted. If an individual is infected with COVID-19, then he/she has to report it to the government. However, there is no obligation to report if a person actually died after the infection. Hence, the data that is obtained is under-reported. In addition, the deaths that are reported may not always be a direct consequence of COVID-19. The person might also have died of other health complications or simply because of an old age. This is seen back in the data set, as some people only die after hundreds of days, which is not likely to be a direct cause of the COVID-19 infection. In order to work with this data set, there have been a few columns added, namely:

– the duration of the illness of the patient, if he/she died.
– the boolean parameter 'death', which is 0 if the person recovered and 1 if he/she died. Note that in this case it is assumed that if death of a person has not been reported, then the person is assumed to have recovered.
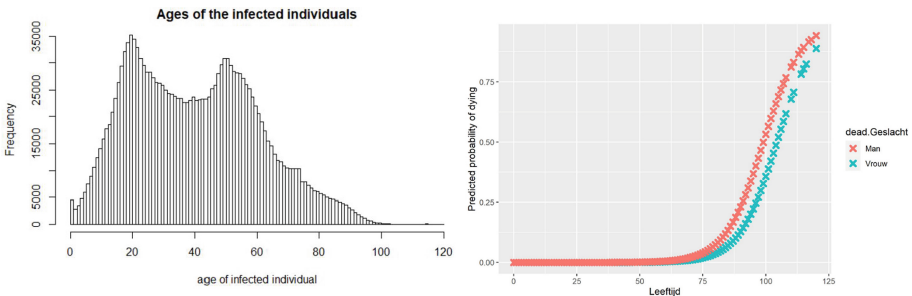


**Fig. 4.** Left: histogram of the age of all the infected individuals. Right: predicted probability of dying against the age of the patient for the genders male, female and undefined

The data ranges from 2020-06-01 until 2021-06-28 and consists of 1666371 observations. The ages range between 0 and 120 years old. There were 98 people who did not report the age. In order to deal with the missing values of the age variable, a random number between 0 and 120 was drawn as the age ranged from 0 till 120. As 98 out of 1666371 observations is a small portion, this missing information did not contaminate the data much. The histogram of the age of the infected individuals is shown in Fig. 4 (Left).

### 5.2  Mathematical Model of the Recovery Rate Estimation

The observations in the data are all assumed to be independent events. In addition, the following variables are defined:

- $l_i$: the age of patient $i$
- $f_i$: the gender of patient $i$, $f_i = 1$ represents female, $f_i = 0$ represents male
- $\delta_i$: $\delta_i = 1$ if the person dies from this COVID-19 infection and $\delta_i = 0$ if the person recovers.
- $x_i$: the duration of the COVID-19 infection after the person got infected, $x_i = y_i$ the duration of the illness and $x_i = \infty$ if the person recovers. So observe that $x_i = \infty \iff \delta_i = 0$. The reason why the duration is set to infinity is because the person recovers from the illness.

The dichotomous variable $\delta_i$ follows a Bernoulli distribution with parameter $p_{l_i,f_i}$, where $p_{l_i,f_i}$ is the probability of dying given the age $l_i$ and gender $f_i$. A logistic regression model is used to estimate this parameter $p_{l_i} = p(l_i, f_i)$, which is actually a function of the age as well. Moreover, it is assumed that given that patient $i$ dies from COVID-19, the duration until death follows an exponential distribution with parameter $\mu_{l_i} = \mu(l_i)$, function of the age. Generalised linear models will be used to estimate this parameter. To begin, we fit a logistic regression model to the data and we apply regression on both age and gender. From the summary we see that both variables, age and gender, are significant as they have very small p-values, suggesting a strong association of the gender as well as the age of a person with the probability of dying. For this case, the following model is obtained for the log odds of death from COVID :

$$\log(\frac{p_{l_i}}{1 - p_{l_i}}) = -13.19 + 0.13l_i - 0.71f_i,$$

where $l_i$ is the age of person $i$ and $f_i = 1$ if the individual is female and $f_i = 0$ if the individual is male. The probability of recovery from COVID-19 is $1 - p(l, f)$. Furthermore, using the fitted model of the log-model, predictions can be made. In Fig. 4 right, predictions of the probability of dying are made based on the age of a person as well as their gender. From these figures it can be seen that males tend to have a higher probability of dying compared to females. In addition, the older the individual, the higher the probability of dying.

### Duration of COVID and Age

Next, we will condition on the fact that the individual has died from COVID, hence all the recovered people are removed from the data frame. In this way the time that it takes for an individual to die, given that the person will die, can be estimated. In the model it is assumed that given that the patient died as a result of this COVID infection, the duration of the illness is exponentially distributed with parameter $\mu_{l_i}$, which depend on the age $l_i$. In this analysis we do not distinguish between gender. The duration of the illness in the data set has the following descriptors:

- Min: 1 day
- 1st Quartile: 7 days
- Median: 11 days
- Mean: 13.71 days
- 3rd Quartile: 16 days
- Max: 391 days.

From these descriptors it can be deduced that individuals who are ill for more than 30 d might not have died as a direct result of the COVID-19 infection. The virus may have caused other health related issues or that the person died because of other reasons than COVID-19. However, in this model it is still assumed that all the individuals who died of COVID-19 have this infection as a cause.

To begin, a linear regression model is fitted to the data. In this case, the assumption is that the observation $x_i$ is drawn from a Normal (Gaussian) distribution with a mean $\mu_i$, that depends on the age $l_i$ and a constant variance $\sigma^2$ across all ages. So $x_i \sim N(\mu_i, \sigma^2)$ and $\mathbb{E}(x_i) = \mu_i = \alpha + \beta l_i$, $\forall i$. Hence the residuals $\epsilon_i = x_i - \mu_i \sim N(0, \sigma^2)$. So actually $x_i = \alpha + \beta l_i + \epsilon_i = 31.70 - 0.22 l_i$, where $\sigma = 12.45$.

Next, a log-transformed linear regression model is applied. This models the duration of the illness on a logarithmic scale, where the model is given by: $\log(x_i) \sim N(\mu_i, \sigma^2)$ and $\mathbb{E}(\log(x_i)) = \mu_i = \alpha + \beta l_i$. The model assumption is that the duration follows a log-normal distribution, so $x_i \sim \log(N(\mu_i, \sigma^2))$ with $\mathbb{E}(x_i) = exp(\mu_i + \sigma^2/2)$.

Finally, a Poisson regression is performed. In this case the response variable $x_i$ is assumed to have a Poisson distribution and it assumes that the $\log(\mathbb{E}(x_i))$ is a linear combination of the unknown parameters. The reason why we use this Poisson regression model is that this model is very common for count data (that is data consisting of natural numbers $\{0, 1, 2, 3, \ldots\}$). In this case, we count the number of days that the individual is ill until the person dies. Although the Poisson model, like the log-transformed, is based on the assumption that the $\log(\mathbb{E}(x_i))$ is a linear combination of the unknown parameters, the main difference with the log-transformed model is that the response variable follows a Poisson distribution, whereas in the log-linear regression model above, the response variable is assumed to follow a Gaussian distribution.

The Poisson distribution only has one parameter $\mu_i$, which is also the expected value. The model is given by: $x_i \sim Poisson(\mu_i)$, $\log(\mu_i) = \alpha + \beta l_i$ and $\mathbb{E}(x_i) = exp(\alpha + \beta l_i)$. The link function is the 'log' function in this case. Note that the mean and the variance are the same for the Poisson probability distribution. The reason why the Poisson distribution is used, is because it will generate integer numbers, which is in line with the actual duration of the data given (even though in reality we do not have whole days of course, but more days + hours + minutes + seconds). In Fig. 5 left, a plot with all the observations, as well as the three different fitted models, is shown. From this plot we see that actually all three models seem to be similar. Interestingly, the data is still quite scattered containing quite a few outliers, which are not detected by any of the models. In order to choose the best model, the Akaike Information

Criterion (AIC) and $R^2$ values of the three models are compared. The values can be found in Table 3, and judging from the AIC and $R^2$–values, it is clear that the log-linear model has the most favourable characteristics.

**Table 3.** Table of the estimated coefficients for the linear, log-linear and Poisson model as well as their null deviance, residual deviance and if present the residual standard deviation. The AIC and $R^2$ values are reported too.

| Model | Linear | Log-linear | Poisson |
|---|---|---|---|
| *Intercept* | 31.70 | 3.31 | 3.81 |
| *Age coefficient* | −0.22 | −0.01 | −0.01 |
| *Null deviance* | 1325070.9 | 4665.2 | 59160.7 |
| *Residual deviance* | 1282943.4 | 4545.8 | 56262.1 |
| *Residual sd sigma* | 12.45 | 0.74 | – |
| *AIC* | 65294.64 | 18544.87 | 91199.97 |
| $R^2$ | 0.0318 | 0.0256 | 0.04890 |

**Statistical Simulations**

Since all the models have a certain distribution associated to it, we subsequently simulate the duration of illness based on the respective three (normal, log-normal and Poisson) probability distributions with the estimated parameters on the actual ages and compare it to the real data. In Fig. 5 right, a plot of the simulated data as well as the real data is shown. From this plot it is clear that the log-transformed linear model performed best (which is consistent with our earlier findings regarding the $R^2$ and AIC statistics) as it is able to also predict higher illness duration (probably due to the variance in the model). Logically, the four outliers of the data are not predicted, but since these are only four points in the data, fitting them into the model as well, would lead to overfitting. Hence from the simulated data it can be seen that a log-linear model is the most suitable model, which is in line with the lowest AIC value.
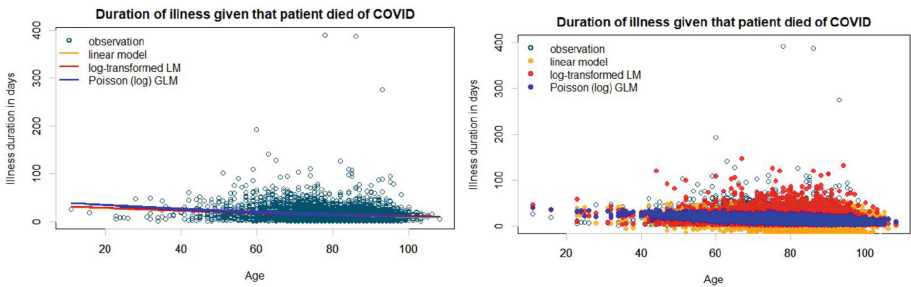


**Fig. 5.** Left: plot of the observations together with the three fitted models: linear, log-linear and Poisson, Right: Simulated duration of the illness duration compared to the actual duration from the data

Therefore, we obtain the following for the duration of illness before death (in days):

$$x_i = \begin{cases} e^{3.31 - 0.012 l_i}, & \text{if the person dies} \\ \infty, & \text{if the person recovers.} \end{cases}$$
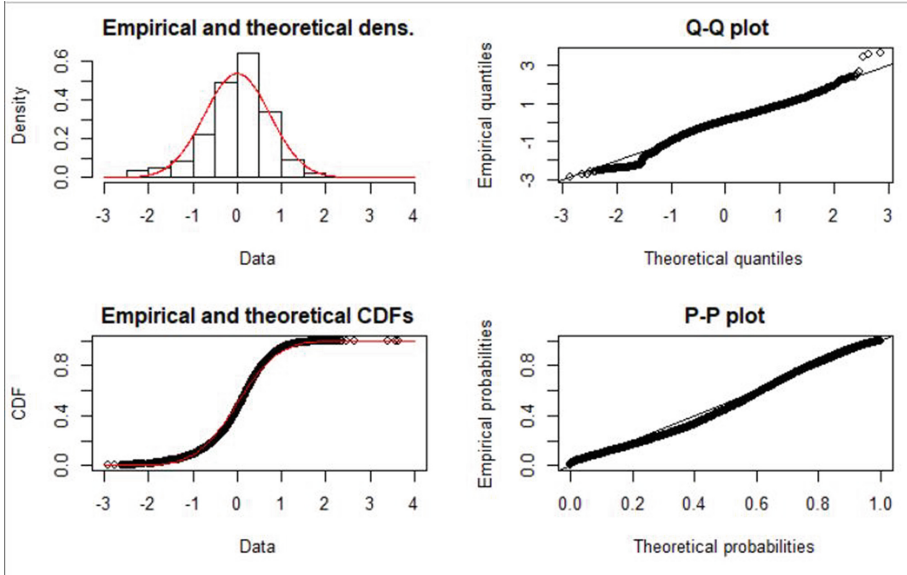


**Fig. 6.** Normal fit of the residuals of the log-linear model. From this we see that the residuals indeed follow a normal distribution.

We have plotted the residuals of the log-linear plot (not shown here), and from this plot it is clear that the residuals are nicely centered around zero. To investigate whether the residuals are normally distributed or not (as it should be according to the assumptions of this model), a normal fit is done and the results are shown in Fig. 6. From this, it is obvious that the residuals are most likely normally distributed. This property indicates that indeed this statistical model is a good fit. We note that the outliers may be disregarded.

## 6    Discussion and Conclusions

We have proposed a mathematical framework to simulate the spread of the novel coronavirus epidemic in a community based on a spatial Markov Chain network model, where the network topology can be made dynamic over time. The model is equipped with uncertainty assessment, in the sense that the probability of different events can be estimated. Different lockdown scenarios have been simulated to see the impact of the severity of lockdown policies. Mild lockdowns were not really effective in reducing the spread of the virus, while heavy lockdowns

caused the number of infected cases to be approximately constant over the lockdown period and severe lockdown protocols could eradicate the virus (under the assumptions of the model). Lifting the lockdown rules caused multiple peaks in the number of active cases as the virus could spread more rapidly when the lockdown rules were relaxed. It is recommended for governments to strictly monitor what is happening when lifting lockdown rules. It might be wise to implement some more stricter rules even if the rules have been relaxed previously, when the spread of the virus is flaming up again. In reality, we see that many governments have indeed taken this approach. The proposed model is different from most of the models that have already been introduced, which are based on the general S(E)IR-models, like in [15–19]. There are two unique features in the presented model: (1) it is based on cellular automata (2) it uses exponentially distributed times between different states, which makes the model stochastic of nature. The model is able to predict the dynamics of the disease under various lockdown scenarios with different infection and recovery rates. However, we also note some limitations of the proposed model. The proposed model is still general as it has not yet been adjusted to specific countries or regions. It is recommended that the simple model is extended and adjusted to specific areas in the world as well as more features regarding the disease are added. This can be done by adding more compartments. To set up the mathematical model, a constant population size of $n$ people is assumed. We assumed that the virus is only transferred from infected people to susceptible people, while in reality the virus is also able to spread through surfaces or objects. People who just recovered from the COVID-19 virus, are able spread the disease too. According to a study published in the journal JAMA, patients who recovered from the COVID-19 virus had been tested positive for the virus in every test between days 5 and 13 post-recovery [20].

The severity of the symptoms caused by the coronavirus vary from person to person. Some people only experience very mild symptoms like fever and dry cough, while others have difficulty breathing, chest pain and might even lose their speech or mobility and must be hospitalised. Elderly (above 70) also have a higher chance of dying from the coronavirus compared to younger individuals. Certain risk groups including individuals with chronic respiratory or pulmonary problems, heart patients or diabetics also have a higher chance of dying as a result of the coronavirus compared to 'normal people' of their age. Hence it is important for further research to investigate the infection rate for the different age-groups as well as the probability of getting heavy or mild symptoms.

Moreover, there is the possibility of reinfection. In this model, the recovery rate $\mu$ is assumed to constant for every individual, while later it is added in the extended model. Research shows that over a period of about three months, the number of antibodies in recovered people rapidly decreases. Recent studies found that there is a high chance of losing immunity to the COVID-19 virus after a period of time [21]. Death is assumed when a person has been infected for a certain period of time. In practice, some people might have mild symptoms for a very long time and recover or have very heavy symptoms and die within a few

days. Variations are large among different individuals. Further research should be carried out to investigate the actual probability of reinfection (which is now only taken as a hypothetical parameter in the simulations) as well as introducing different sub-groups of infected people based on their symptoms, age and other health conditions (like asthma). To better quantify these transmission rates, one could assume that they are functions of time, since the rate of recovery does not only change with age, but also with the circumstances at that time (for instance IC capacity).

The different lockdown simulations showed that a severe lockdown is able to extinguish the outbreak in a limited amount of time, while less severe lockdown policies mainly cause a steady amount of infections as well as recovering and dying people over time. The term *flattening the curve* is often used in the media to describe: (i) reduction in the peak number of infections, to prevent the health care system from being overloaded and (ii) increasing the duration of the pandemic over a larger time interval, but with the same number of cases at the end. This phenomenon has been seen back in the simulations for various lockdown scenarios. The time of the epidemic was stretched over a longer period when a lockdown was implemented and the peak number of infections decreased. However, the total number of infections was the same as well as the total number of people who recovered and died, as can be seen in the cumulative graphs. Lifting the lockdown rules resulted into several peaks in the number of infections. The second and possibly third peak were a lot lower compared to the first one, due to a lower number of susceptible people. Nonetheless, governments are encouraged to impose appropriate measures in their lockdown policies to reduce the peak in the number of infections when they are relaxing the lockdown rules.

The model provides a theoretical framework to investigate the spread of the COVID-19 virus. The variability in the results due to the randomness in the model, makes the simulations more realistic as similar lockdown protocols in different countries have different effects on the number of casualties. Culture, population density and public health care are examples of variables that have a major impact on the number of casualties. In the paper by Cooper, Mondal and Antonopoulos in [22] a SIR-model is developed for various different communities. The paper by Cooper et al. present a study of the time evolution of different populations and the diversity in the parameters for the spread of the disease. Future research on this topic could help modify this model to a specific country or region. We realise that the presented model is still very basic and more simulations need to be carried out using extensions of the basic model. Perhaps it will be possible to find a pair of values for $(\lambda_g, \mu)$ and a lockdown strategy to eradicate the virus. It is also advisable to consider larger population samples or more simulation runs, since in this case we have only performed simulations for a constant population size of 10000 people and we have only done 100 simulation runs per pair of $(\lambda_g, \mu)$.

As the presented model is not yet adjusted to a specific country or region neither is it as extensive as many proposed S(E)IR-models. The results cannot be directly compared to previous studies. In future studies, it might be possible

to relate the outcomes of this Spatial Cellular Automata Markov Chain model to the S(E)IR-models presented in for instance [15–19]. A downside of the current model, compared to the classical deterministic models like the S(E)IR-models, is that it is relatively expensive to execute. One could possibly optimise the model and use more computational resources within the Monte Carlo framework.

Estimation of the infection rate parameter for each individual is difficult based on the observed data. This is because the probability of getting infected by COVID-19 is dependent on many factors including the number of contacts that a person has (which in simulations must be randomised), the incubation period, lockdown policies, vaccination and many other conditions. In addition, the data is subject to under-reporting, due to the fact that in the beginning of the pandemic there were not that many tests performed and not everybody who becomes infected with COVID-19 reports it. This can be due to having had mild or almost no symptoms. The effect of lockdown is seen back in the number of active cases and in the Netherlands almost every two weeks the lockdown policy has changed in the first year of the pandemic. In addition, the number of tests increased over time, which resulted into more people who were tested positive for the virus compared to earlier. The rise of the different variants of COVID-19 as well as the increase in vaccinated people need to be taken in consideration as well. Further studies, which take these variables into account, will need to be carried out.

Currently, almost all the parameter estimation has been done based on statistical regression models such as the logistic regression, linear regression, log-linear regression and log-Poisson regression for the probability of recovery and duration of the disease given that the patient dies as a function of age and gender. In order to have a direct correspondence between the Markov Chain model and the data, more research is needed. For instance in [23] a physics-informed neural network is used to estimate the time-dependent contact rate and in [24] a Bayesian framework is used to estimate the infection parameter. Future research is encouraged to try to use any statistical method to estimate the parameter based on the mathematical model presented in this paper. Additionally, the estimation of the recovery rate parameter $\mu$ based on the actual data of the RIVM is also biased due to the assumptions made in the model. Further research should be undertaken to investigate this parameter under preferably less strong assumptions. All together, the proposed mathematical model is different from the regular S(E)IR-models used by many countries at this moment. It might provide a different way of modelling the pandemic and potentially lead to more accurate or different predictions as it looks at each individual individually. However, the proposed model is still simplistic and has to be extended and applied to real data in order to make predictions and give conclusions on the evolution of this disease.

# References

1. worldometer. COVID-19 CORONAVIRUS PANDEMIC. https://www.worldometers.info/coronavirus/?utm_campaign=homeAdvegas1?%22%20%5Cl%20%22countries

2. WHO. Naming the coronavirus disease (COVID-19) and the virus that causes it (2020). https://www.who.int/emergencies/diseases/novel-coronavirus-2019/technical-guidance/naming-the-coronavirus-disease-(covid-2019)-and-the-virus-that-causes-it

3. Staff, M.C.: Covid-19 (coronavirus): Long-term effects, May 2021. https://www.mayoclinic.org/diseases-conditions/coronavirus/in-depth/coronavirus-long-term-effects/art-20490351

4. Kermack, W.O., McKendrick, A.G.: A contribution to the mathematical theory of epidemics. In: Proceedings of the royal society of London. Series A, Containing Papers of a Mathematical and Physical Character, vol. 115, no. 772, pp. 700–721 (1927)

5. Getz, W.M., Salter, R., Muellerklein, O., Yoon, H.S., Tallam, K.: Modeling epidemics: a primer and numerus model builder implementation. Epidemics **25**, 9–19 (2018)

6. Alleman, T.W., Vergeynst, J., De Visscher, L., Rollier, M., Torfs, E., Nopens, I., Baetens, J.: Assessing the effects of non-pharmaceutical interventions on SARS-CoV-2 transmission in Belgium by means of an extended SEIQRD model and public mobility data. Epidemics **37**, 100505 (2021)

7. Allen, L.J.: A primer on stochastic epidemic models: formulation, numerical simulation, and analysis. Infect. Dis. Model. **2**(2), 128–142 (2017)

8. Pellis, L., et al.: Eight challenges for network epidemic models. Epidemics **10**, 58–62 (2015)

9. Walters, C.E., Meslé, M.M., Hall, I.M.: Modelling the global spread of diseases: a review of current practice and capability. Epidemics **25**, 1–8 (2018)

10. Duan, W., Fan, Z., Zhang, P., Guo, G., Qiu, X.: Mathematical and computational approaches to epidemic modeling: a comprehensive review. Front. Comput. Sci. **9**(5), 806–826 (2014). https://doi.org/10.1007/s11704-014-3369-2

11. Britton, T.: Stochastic epidemic models: a survey. Math. Biosci. **225**(1), 24–35 (2010)

12. O'Neill, P.D.: A tutorial introduction to Bayesian inference for stochastic epidemic models using Markov chain monte Carlo methods. Math. Biosci. **180**(1–2), 103–114 (2002)

13. Vermolen, F., Pölönen, I.: Uncertainty quantification on a spatial markov-chain model for the progression of skin cancer. J. Math. Biol. **80**(3), 545–573 (2020)

14. Chen, F.: Better modelling of infectious diseases: lessons from covid-19 in china. BMJ **375**, 2363 (2021)

15. Fanelli, D., Piazza, F.: Analysis and forecast of covid-19 spreading in china, Italy and France. Chaos, Solitons Fractals **134**, 109761 (2020)

16. Yang, C., Wang, J.: A mathematical model for the novel coronavirus epidemic in Wuhan, china. Math. Biosci. Eng. **17**(3), 2708–2724 (2020)

17. Caccavo, D.: Chinese and italian covid-19 outbreaks can be correctly described by a modified sird model, *medRxiv* (2020)

18. Al-Raeei, M.: The forecasting of covid-19 with mortality using SIRD epidemic model for the united states, Russia, China, and the Syrian Arab republic. AIP Adv. **10**(6), 065325 (2020)

19. Rajagopal, K., Hasanzadeh, N., Parastesh, F., Hamarash, I.I., Jafari, S., Hussain, I.: A fractional-order model for the novel coronavirus (covid-19) outbreak. Nonlinear Dyn. **101**(1), 711–718 (2020)

20. Lan, L., et al.: Positive RT-PCR test results in patients recovered from COVID-19. JAMA **323**(15), 1502–1503 (2020). https://doi.org/10.1001/jama.2020.2783

21. Agel, F.: Antibodies, immunity low after COVID-19 recovery. https://www.dw.com/en/coronavirus-antibodies-immunity/a-54159332
22. Cooper, I., Mondal, A., Antonopoulos, C.G.: A sir model assumption for the spread of covid-19 in different communities. Chaos, Solitons Fractals **139**, 110057 (2020)
23. Grimm, V., Heinlein, A., Klawonn, A., Lanser, M., Weber, J.: Estimating the time-dependent contact rate of sir and seir models in mathematical epidemiology using physics-informed neural networks," Universität zu Köln, Technical Report, September 2020. https://kups.ub.uni-koeln.de/12159/
24. Irons, N.J., Raftery, A.E.: Estimating sars-cov-2 infections from deaths, confirmed cases, tests, and random surveys. In: Proceedings of the National Academy of Sciences, vol. 118, no. 31 (2021). https://www.pnas.org/content/118/31/e2103272118