

GreatMod: SIR model

Beccuti Marco, Castagno Paolo, Pernice Simone

Contents

Introduction	1
How to start	2
Something to know	2
Cases of study	2
SIR model	3
Analysis	3
Model generation	3
Model Analysis	4
Calibration analysis	7
References	12

Introduction

In this document we describe how to use the R library *epimod*. In details, *epimod* implements a new general modeling framework to study epidemiological systems, whose novelties and strengths are:

1. the use of a graphical formalism to simplify the model creation phase;
2. the automatic generation of the deterministic and stochastic process underlying the system under study;
3. the implementation of an R package providing a friendly interface to access the analysis techniques implemented in the framework;
4. a high level of portability and reproducibility granted by the containerization (Veiga Leprevost et al. 2017) of all analysis techniques implemented in the framework;
5. a well-defined schema and related infrastructure to allow users to easily integrate their own analysis workflow in the framework.

The effectiveness of this framework is showed through the wellknown and simple SIR model.

How to start

Before starting the analysis we have to install (1) GreatSPN GUI, (2) docker, and (3) the R package **devtools** for installing *EPIMOD*. GreatSPN GUI, the graphical editor for drawing Petri Nets formalisms, is available online ([link to install GreatSPN](#)), and it can be installed following the steps showed therein. Then, the user must have docker installed on its computer for exploiting the *epimod*'s docker images (for more information on the docker installation see: [link to install docker](#)), and to have authorization to execute docker commands reported in the command page of function `install docker`. To do this the following commands must be executed.

1. Create the docker group.

```
$ sudo groupadd docker
```

2. Add your user to the docker group.

```
$ sudo usermod -aG docker $USER
```

The R package *devtools* has to be installed to run *epimod*:

```
install.packages("devtools")  
library(devtools)  
install_github("qBioTurin/epimod", dependencies=TRUE)
```

```
library(epimod)
```

Then, the following function must be used to download all the docker images used by *epimod*:

```
downloadContainers()
```

Something to know

All the *epimod* functions print the following information:

- *Docker ID*, that is the CONTAINER ID which is executed by the function;
- *Docker exit status*, if 0 then the execution completed with success, otherwise an error log file is saved in the working directory.

Cases of study

In this section we show the steps necessary to model, study and analyze a simple case study. To this aim, we choose to study the diffusion of a disease following the SIR dynamics. We refer to (Keeling and Rohani 2011) for all the details.

SIR model

The S-I-R models the diffusion of an infection in a population assuming three possible states (or compartments) in which any individual in the population may move. Specifically, (1) *Susceptible*, individuals unexposed to the disease, (2) *Infected*, individuals currently infected by the disease, and (3) *Recovered*, individuals which were successfully recovered by the infection. To consider the simplest case, we ignore the population demography (i.e., births and deaths of individuals are omitted), thus we consider only two possible events: the infection (passage from *Susceptible* to *Infected*), and the recovery (passage from *Infected* to *Recovered*). We are also assuming to neglect complex pattern of contacts, by considering an homogeneous mixing. From a mathematical point of view, the system behaviors can be investigated by exploiting the deterministic approach (Kurtz 1970) which approximates its dynamics through a system of ordinary differential equations (ODEs):

$$\begin{aligned}\frac{dS}{dt} &= -\beta SI, \\ \frac{dI}{dt} &= \beta SI - \gamma I, \\ \frac{dR}{dt} &= \gamma I,\end{aligned}$$

where:

- S , I , R are the number of susceptible, infected, and recovered individuals, respectively;
- β is the infection rate;
- γ is the recovery rate, which determines the mean infectious period.

Analysis

Model generation

The first step is the model construction. Starting with the GreatSPN editor tool it is possible to draw the model using the PN formalism and its generalizations. We recall that the Petri Nets are bipartite graphs in which we have two type of nodes, places and transitions. Graphically, places are represented as circles and those are the variables of our systems. On the other hand, transitions are depicted as rectangles and are the possible events happening in the system. Variables and events (i.e., places and transitions) are connected through arcs, showing what variable(s) is (are) affected by a specific event. For more details we refer to (Marsan et al. 1995).

Therefore, as represented in figure 1, we add one place for each variable of the system (i.e., S , I , and R represent the susceptible, infected, and recovered individuals respectively), and one transition for each possible event (i.e., *Infection* and *Recovery*). Finally, we save the PN model as a file with extension *.PNPRO*.

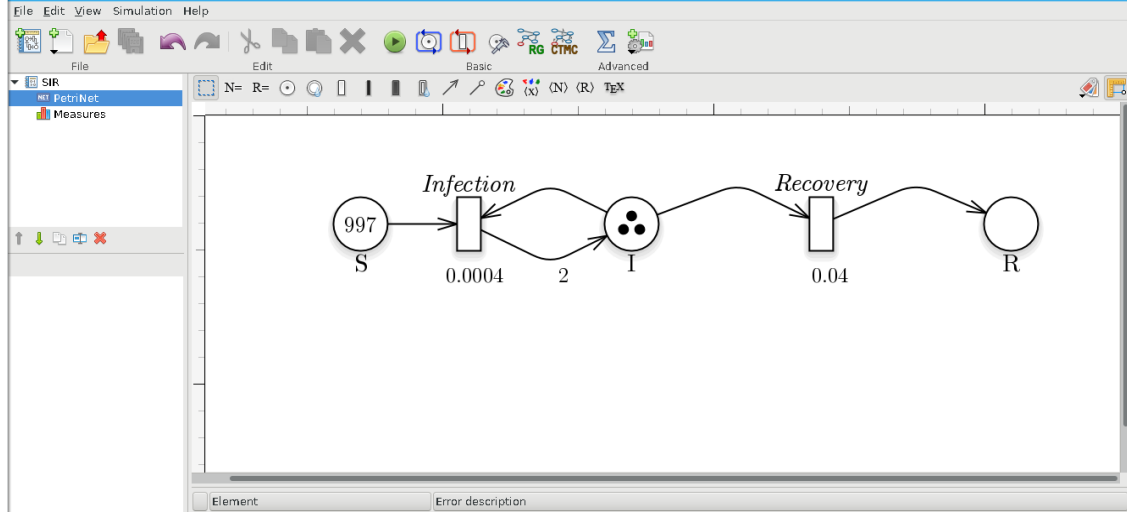


Figure 1: Petri Net representation of the SIR model.

Having constructed the model, the generation of both the stochastic (the Continuous Time Markov Chain) and deterministic (ODEs) processes underlying the model is implemented by the *model_generation()* function. This function takes as input the file generated by the graphical editor, in this case called *SIR.PNPRO*, and automatically derives the processes.

```
library(epimod)
model_generation("SIR.PNPRO")
```

The binary file *SIR.solver* is generated in which the derived processes and the library used for their simulation are packaged.

Model Analysis

A possible step is the model analysis, where the corresponding function *model_analysis()* executes and tests the behavior of the developed model. Furthermore, by changing the input parameters, it is possible to perform a *what-if* analysis or forecasting the evolution of the diffusion process.

```
model_analysis(solver_fname = "SIR.solver",
               parameters_fname = "Functions_list_ModelAnalysis.csv",
               f_time = 100, # days
               s_time = 1,
               )
```

1. **solver_fname**: the *.solver* file generated by the *model_generation* function, that is *SIR.solver*;
2. **parameters_fname**: a textual file in which the parameters to be studied are listed associated with their range of variability. This file is defined by three mandatory columns: (1) a tag representing the parameter type: *i* for the complete initial marking (or condition), *p* for a single parameter (either a single rate or initial marking), and *g* for a rate associated with general transitions (Pernice et al. 2019) (the user must define a file name coherently with the one used in the general transitions file); (2) the name of the transition which is varying (this must correspond to name used in the PN draw in GreatSPN editor), if the complete initial marking is considered (i.e., with tag *i*) then by default the name *init* is used; (3) the function used for sampling the value of the variable considered, it could be

either a R function or an user-defined function (in this case it has to be implemented into the R script passed through the *functions_fname* input parameter). Let us note that the output of this function must have size equal to the length of the varying parameter, that is 1 when tags *p* or *g* are used, and the size of the marking (number of places) when *i* is used. The remaining columns represent the input parameters needed by the functions defined in the third column.

1	p;Recovery;0.04;
2	p;Infection;0.0004;
3	

3. **f_time**: the final solution time, for instance 100 days;
4. **s_time**: the time step defining the frequency at which explicit estimates for the system values are desired, in this case it could be set to 1 day.

This function solves the system given a specific parameters configuration which is passed through the function parameter, *parameters_fname*.

```
## How to generate the plots

source("Rfunction/ModelAnalysisPlot.R")

AnalysisPlot = ModelAnalysisPlot(solverName_path = "./results_model_analysis/SIR-analysys-1.trace",
                                Stoch = F,
                                print=F)

AnalysisPlot$plAll
```

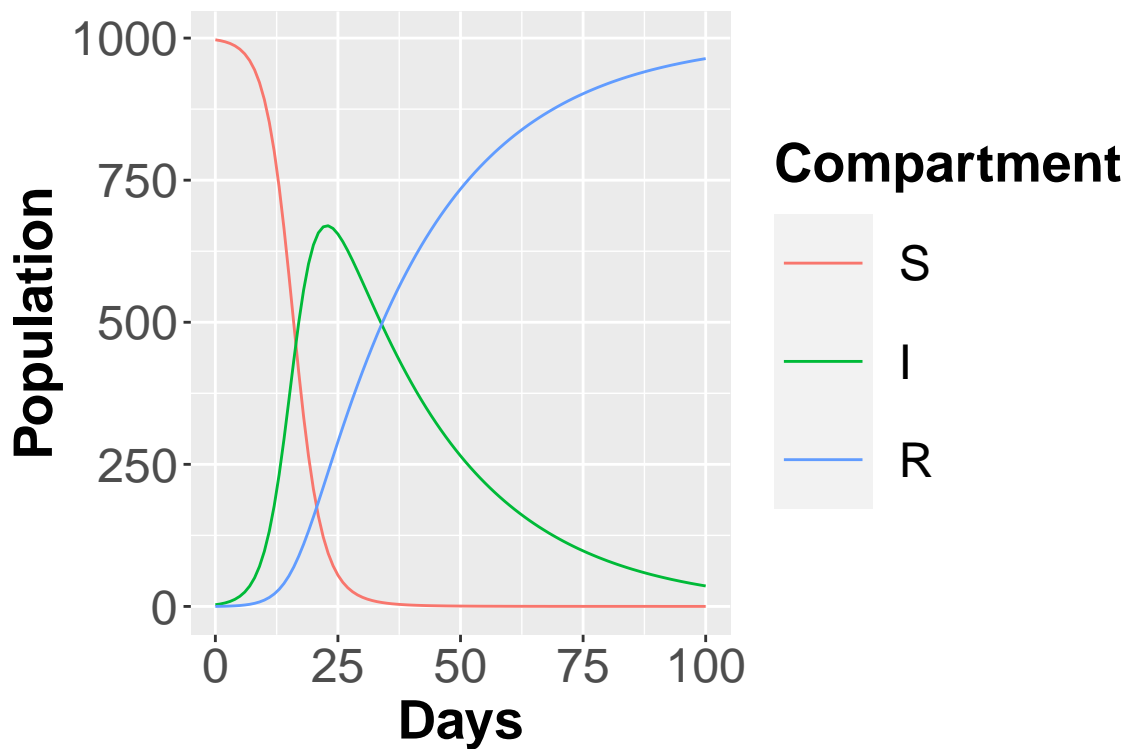


Figure 2: Deterministic Trajectory considering all places

It is possible also to simulate the stochastic behavior of the system exploiting the Gillespie algorithm, namely **SSA**, which is an exact stochastic method widely used to simulate chemical systems whose behavior can be described by the Master equations.

```
model_analysis(solver_fname = "SIR.solver",
               parameters_fname = "Functions_list_ModelAnalysis.csv",
               solver_type = "SSA",
               n_run = 500,
               parallel_processors = 2,
               f_time = 100, # days
               s_time = 1
            )
```

1. **solver_type**: type of solver to use;
2. **n_run**: number of stochastic simulations to run.

How to generate the plots

```
source("Rfunction/ModelAnalysisPlot.R")

AnalysisPlot = ModelAnalysisPlot(solverName_path = "./results_model_analysis/SIR-analysis-1.trace",
                                Stoch = T)

AnalysisPlot$plAll
AnalysisPlot$plAllMean
```

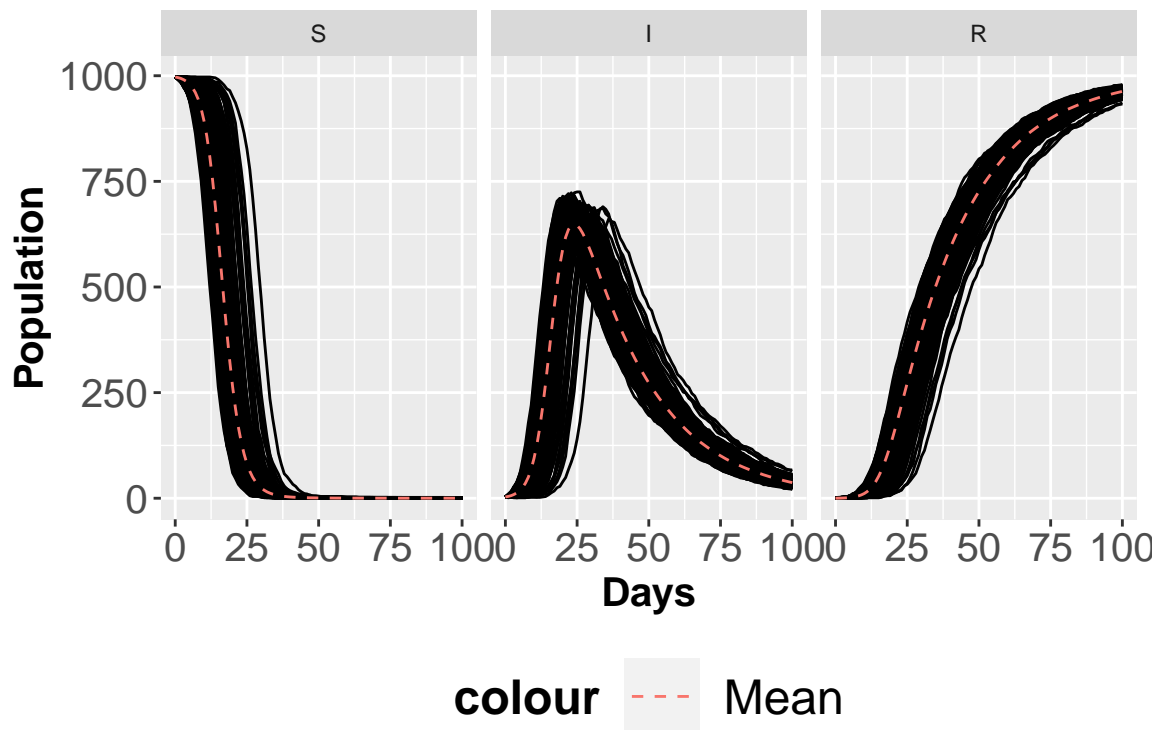


Figure 3: Stochastic Trajectories considering the S place.

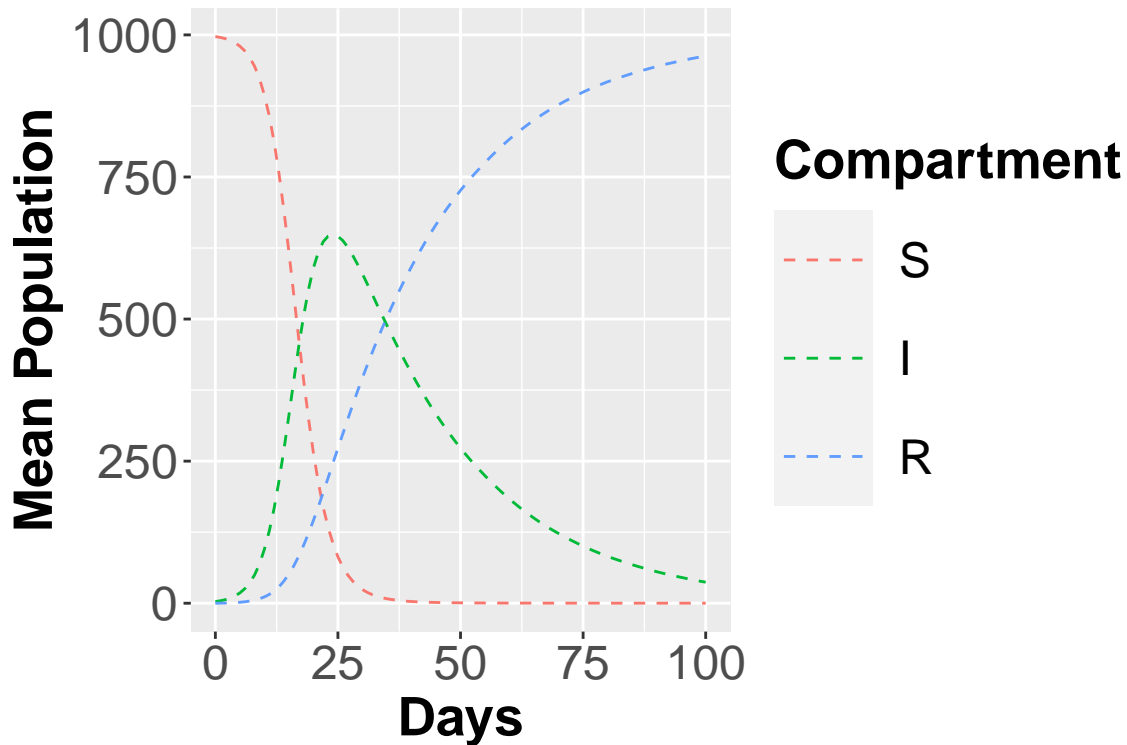


Figure 4: Stochastic Trajectories considering the I place.

What-if analysis It is possible to change the place marking directly from the *parameters_fname*:

```

1 p;I;10;
2 p;S;990;
3 p;Recovery;0.04;
4 p;Infection;0.0004;
5 |

```

Calibration analysis

The aim of this phase is to optimize the fit of the simulated behavior to the reference data by adjusting the parameters associated with both Recovery and Infection transitions. This step is performed by the function *model_calibration()*, characterized by the solution of an optimization problem in which the distance between the simulated data and the reference data is minimized, according to the definition of distance provided by the user (*distance_fname*).

```

model_calibration(solver_fname = "SIR.solver",
                  parameters_fname = "Functions_list_Calibration.csv",
                  functions_fname = "FunctionCalibration.R",
                  reference_data = "reference_data.csv",
                  distance_measure_fname = "mse.R" ,
                  f_time = 100, # days
                  s_time = 1, # day
                  # Vectors to control the optimization

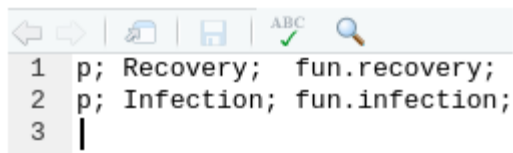
```

```

ini_v = c(0.035,0.00035),
ub_v = c(0.05, 0.0005),
lb_v = c(0.025, 0.00025),
max.time = 1
)

```

1. **solver_fname**: the *.solver* file generated by the *model_generation* function, that is *SIR.solver*;
2. **parameters_fname**: a textual file in which the parameters to be studied are listed associated with their range of variability. This file is defined by three mandatory columns (*which must be separated using ;*): (1) a tag representing the parameter type: *i* for the complete initial marking (or condition), *p* for a single parameter (either a single rate or initial marking), and *g* for a rate associated with general transitions (Pernice et al. 2019) (the user must define a file name coherently with the one used in the general transitions file); (2) the name of the transition which is varying (this must correspond to name used in the PN draw in GreatSPN editor), if the complete initial marking is considered (i.e., with tag *i*) then by default the name *init* is used; (3) the function used for sampling the value of the variable considered, it could be either a R function or an user-defined function (in this case it has to be implemented into the R script passed through the *functions_fname* input parameter). Let us note that the output of this function must have size equal to the length of the varying parameter, that is 1 when tags *p* or *g* are used, and the size of the marking (number of places) when *i* is used. The remaining columns represent the input parameters needed by the functions defined in the third column. An example is given by the file *Functions_list_Calibration.csv*:



```

1 p; Recovery; fun.recovery;
2 p; Infection; fun.infection;
3 |

```

where the rates of the *Recovery* and *Infection* transitions can be calibrated by using the R functions stored in the R script *functions_fname*;

3. **functions_fname**: an R file storing the user defined functions to generate instances of the parameters summarized in the *parameters_fname* file. An example is given by *FunctionCalibration.R*, where the function *recovery.fun* introduced in *Functions_list_Calibration.csv* file is defined in order to return the value (or a linear transformation) of the vector of the unknown parameters generated from the optimization algorithm, namely \mathbf{x} , whose size is equal to number of parameters in *parameters_fname*. Let us note that the output of these functions must return a value for each input parameter. For instance, to calibrate the transition rates associated with *Recovery* and *Infection*, the functions *recovery* and *infection* have to be defined, returning just the corresponding value from the vector \mathbf{x} , where $\mathbf{x}[1]$ = “Recovery rate”, $\mathbf{x}[2]$ = “Infection rate”, *since we do not want to change the vector generated from the optimization algorithm. The order of values in \mathbf{x} is given by the order of the parameters in parameters_fname*”.

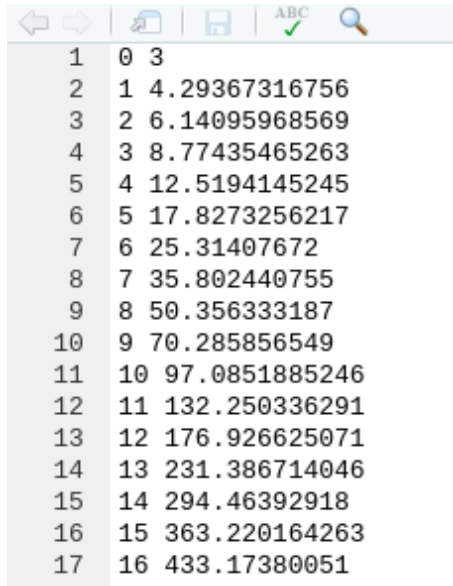
```

fun.recovery<-function(x)
{
  return(x[1])
}

fun.infection<-function(x)
{
  return(x[2])
}

```


4. **reference_data**: a csv file storing the data to be compared with the simulations' result. In *reference_data.csv* we report the SIR evolution starting with 997 susceptible, three infected and zero recovered, with a recovery and infection rates equals to 0.04 and 0.0004 respectively. Notice that the **reference_data**'s rows must be the variable time serie, and so the columns the corresponding values at a specific time.



1	0	3
2	1	4.29367316756
3	2	6.14095968569
4	3	8.77435465263
5	4	12.5194145245
6	5	17.8273256217
7	6	25.31407672
8	7	35.802440755
9	8	50.356333187
10	9	70.285856549
11	10	97.0851885246
12	11	132.250336291
13	12	176.926625071
14	13	231.386714046
15	14	294.46392918
16	15	363.220164263
17	16	433.17380051

5. **distance_measure_fname**: the R file storing the function to compute the distance (or error) between the model output and the reference dataset itself. The function defining the distance takes in input only the reference data (defined in *reference_data.csv*), and the *simulation output* with the following structure:

	Time	S	I	R
1	0	997.0000	3.000000	0.000000
2	1	995.5028	4.356570	0.1405834
3	2	993.3335	6.321825	0.3446604
4	3	990.1957	9.163666	0.6406362
5	4	985.6685	13.262202	1.0693255
6	5	979.1604	19.150585	1.6890500
7	6	969.8535	27.564066	2.5824805

an example is given by *mse.R* where a distance measure (based on the squared error distance) as function of the infected individuals is defined:

```
mse<-function(reference, output)
{
  reference[1,] -> times_ref
  reference[2,] -> infect_ref

  # We will consider the same time points
  Infect <- output[which(output$Time %in% times_ref),"I"]
  infect_ref <- infect_ref[which( times_ref %in% output$Time)]

  diff.Infect <- 1/length(times_ref)*sum(( Infect - infect_ref )^2 )
}
```

```

    return(diff.Infect)
}

```

6. **f_time**: the final solution time, for instance 10 weeks (70 days);
7. **s_time**: the time step defining the frequency at which explicit estimates for the system values are desired, in this case it could be set to 1 day;
8. **ini_v**: Initial values for the parameters to be optimized.
9. **lb_v, ub_v**: Vectors with length equal to the number of parameters which are varying. Lower/Upper bounds for each parameter.
10. **max_time**: maximum running time.

Let us observe that: the distance and target functions must have the same name of the corresponding R file.

```

# How to generate the plots

source("Rfunction/CalibrationPlot.R")

plots <- calibration.plot(solverName_path = "./results_model_calibration/SIR-calibration-1.trace",
                          reference_path = "reference_data.csv")

plots$plS
plots$plI
plots$plR

```

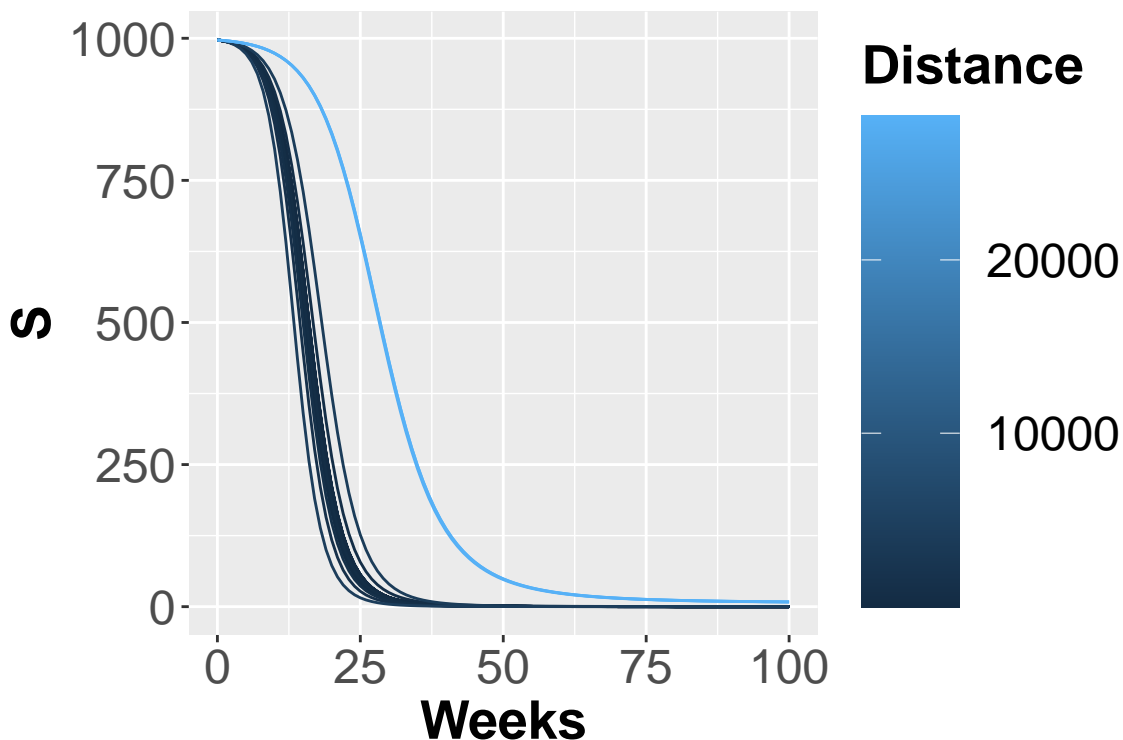


Figure 5: Trajectories considering the S place.

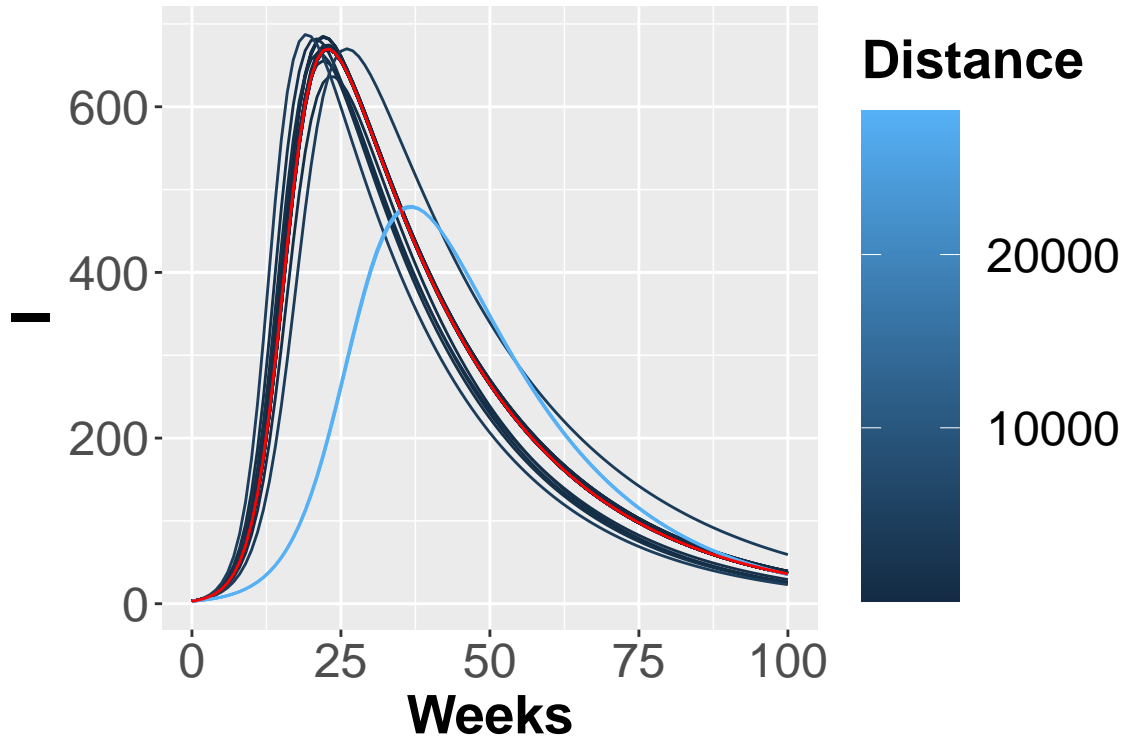


Figure 6: Trajectories considering the I place.

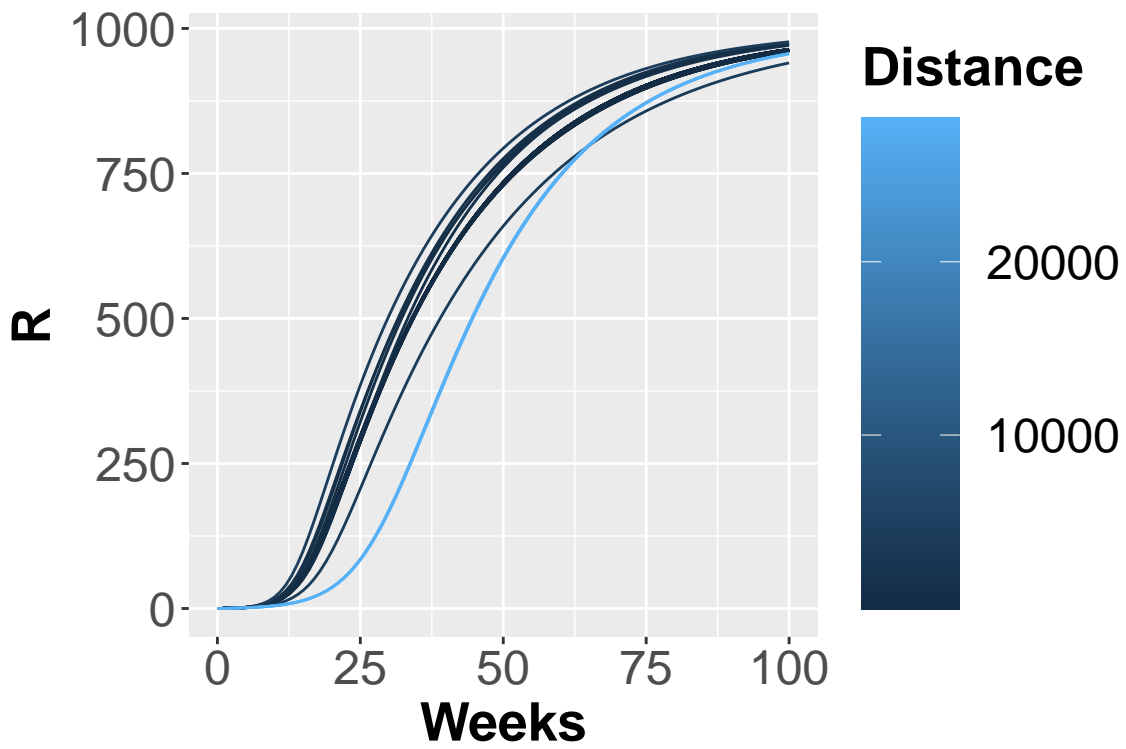


Figure 7: Trajectories considering the R place.

In figures 6,5 and 7 the trajectories with color depending on the squared error w.r.t. reference trend are

plotted. In this case, fixing a maximum number of objective function calls, we obtain the following optimal value for the two parameters:

[1] 0.0399999567 0.0004000018

References

- Keeling, Matt J, and Pejman Rohani. 2011. *Modeling Infectious Diseases in Humans and Animals*. Princeton University Press.
- Kurtz, T. G. 1970. "Solutions of Ordinary Differential Equations as Limits of Pure Jump Markov Processes." *J. Appl. Probab.* 1 (7): 49–58.
- Marsan, M. Ajmone, G. Balbo, G. Conte, S. Donatelli, and G. Franceschinis. 1995. *Modelling with Generalized Stochastic Petri Nets*. New York, NY, USA: J. Wiley.
- Pernice, S., M. Pennisi, G. Romano, A. Maglione, S. Cutrupi, F. Pappalardo, G. Balbo, M. Beccuti, F. Cordero, and R. A. Calogero. 2019. "A Computational Approach Based on the Colored Petri Net Formalism for Studying Multiple Sclerosis." *BMC Bioinformatics*.
- Veiga Leprevost, Felipe da, Björn A Grüning, Saulo Alves Afitos, Hannes L Röst, Julian Uszkoreit, Harald Barsnes, Marc Vaudel, et al. 2017. "BioContainers: an open-source and community-driven framework for software standardization." *Bioinformatics* 33 (16): 2580–2.