

# **R4PDE**

Emerson M. Del Ponte

4/22/2022

# Table of contents

<b>About</b>	<b>3</b>
<b>1 Temporal analysis</b>	<b>4</b>
1.1 Introduction . . . . .	4
1.1.1 Population dynamics models . . . . .	6
1.2 Non-flexible models . . . . .	7
1.2.1 Exponential . . . . .	7
1.2.2 Monomolecular . . . . .	10
1.2.3 Logistic . . . . .	12
1.2.4 Gompertz . . . . .	15
1.3 Model fitting . . . . .	17
1.3.1 Non-replicated . . . . .	18
1.3.2 Designed experiments . . . . .	26
<b>2 Spatial gradients</b>	<b>37</b>
2.1 Models . . . . .	37
2.1.1 Exponential model . . . . .	37
2.1.2 Power law model . . . . .	39
2.2 Linearization of the models . . . . .	41
2.2.1 Transformations of y . . . . .	41
2.2.2 Plot for the linearized form of models . . . . .	42
2.3 Model fitting . . . . .	45
2.3.1 Dataset . . . . .	45
2.3.2 Visualize the gradient . . . . .	46
2.3.3 Linear regression . . . . .	47
<b>3 Spatial patterns</b>	<b>53</b>
3.1 Introduction . . . . .	53
3.2 Position and status within a row (plant to plant) . . . . .	54
3.2.1 Runs test . . . . .	54
3.2.2 Doublets . . . . .	57
3.2.3 Join count . . . . .	58
<b>References</b>	<b>63</b>

# About

**R for Plant Disease Epidemiology** (R4PDE) is a book project in the development stage. It is based on my teaching notes of a graduate course, currently **FIP 602 - Plant Disease Epidemiology**, offered every year for students of the [Graduate Program in Plant Pathology](#) of the Universidade Federal de Viçosa.

This book is for those interested in studying and modelling plant disease epidemics using R programming. Here, I provide context and showcase several methods for describing, visualizing and analyzing epidemic data collected over time and/or space.

Users should have a minimum knowledge of R programming using the RStudio IDE and the tidyverse ecosystem. I make use of custom functions or general and specific R packages for conducting the most common tasks related with the analysis of plant disease epidemiology data.

# 1 Temporal analysis

## 1.1 Introduction

A key understanding of the epidemics relates to the knowledge of rates and patterns. Epidemics can be viewed as dynamic systems that change their state as time goes. The first and simplest way to characterize such changes in time is to produce a graphical plot called disease progress curve (DPC). This curve can be obtained as long as the intensity of the disease ( $y$ ) in the host population is assessed sequentially in time ( $t$ ).

A DPC summarizes the interaction of the three main components of the disease triangle occurring during the epidemic. The curves can vary greatly in shape according to variations in each of the components, in particular due to management practices that alter the course of the epidemics and for which the goal is to stop disease increase. We can create a dataframe in R for a single DPC and make a plot using ggplot. By convention we use  $t$  for time and  $y$  for disease intensity.

Firstly, let's load the essential R packages and set up the environment.

```
library(tidyverse) # loads the essential packages
library(patchwork) # for multipanel figures
theme_set(theme_light()) # set global theme
```

There are several ways to create a dataframe in R. I like to use the `tribble` function as below. The entered data will be assigned to a dataframe called `dpc`.

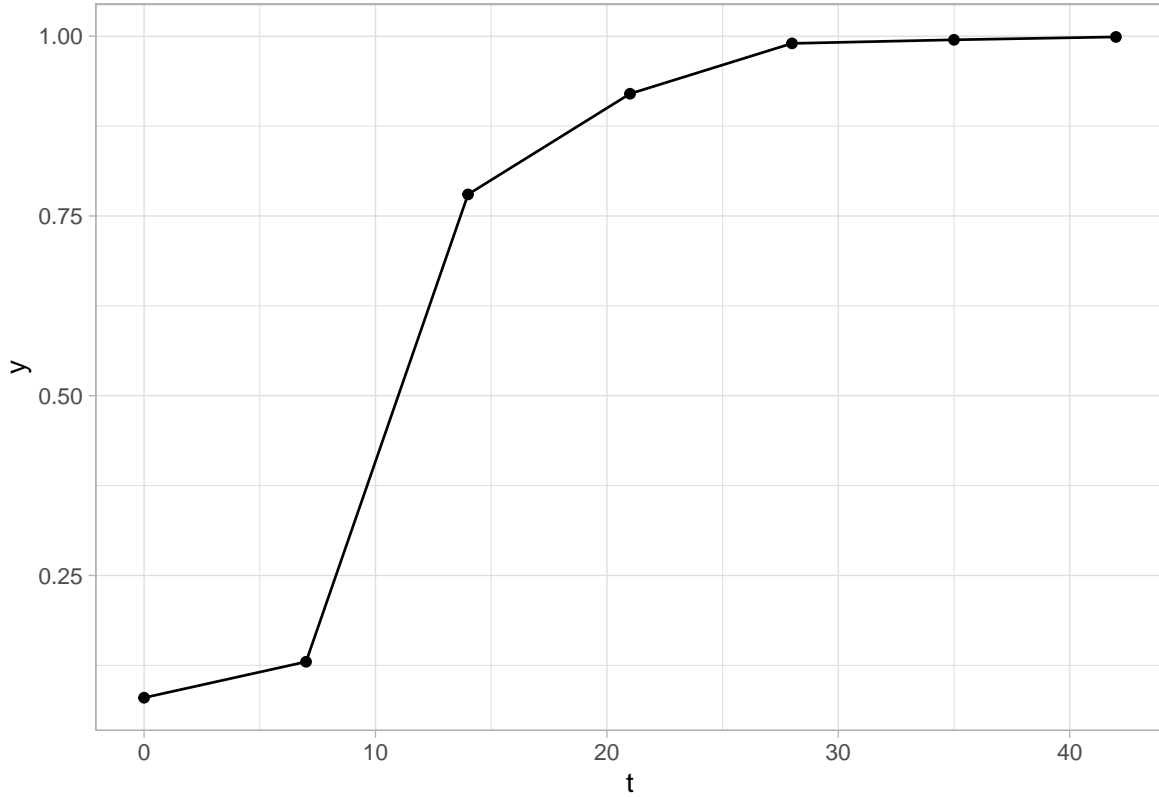
```
dpc <-  
  tribble(  
    ~t, ~`y`,  
    0, 0.08,  
    7, 0.13,  
    14, 0.78,  
    21, 0.92,  
    28, 0.99,  
    35, 0.995,  
    42, 0.999,  
  )
```

Now the plot

**i** Note

Note that I use [pipes](#) in my programming to express a sequence of multiple operations in a more intuitive way.

```
dpc %>%  
  ggplot(aes(t, y)) +  
  geom_point()+  
  geom_line()
```



The depiction and analysis of disease progress curves can provide useful information for gaining understanding of the underlying epidemic process. In a practical manner, the curves are extensively used to evaluate how disease control measures affect epidemics. When characterizing DPCs, a researcher may be interested in describing and comparing epidemics that result from different treatments, or simply in their variations as affected by changes in environment, host or pathogen.

Mathematical models can be fitted to the DPC data to express epidemic progress in terms of rates and absolute/relative quantities. The latter can be accomplished using population dynamics (or growth-curve) models for which the estimated parameters are usually meaningful biologically and appropriately describe epidemics that do not decrease in disease intensity. By fitting an appropriate model to the progress curve data, another set of parameters is available to the researcher when attempting to represent, understand or compare epidemics.

### 1.1.1 Population dynamics models

The family of models that describe the growth of epidemics, hence population dynamics model, are known as deterministic models of continuous time (Madden et al. 2017a). These models

are usually fitted to DPC data to obtain two or more biologically meaningful parameters.

In this tutorial, these models and their formulations are shown using R scripts to simulate the theoretical curves for each model. The reader should be capable of opening R or RStudio software and reproduce the analysis by copying and pasting the codes. Hence, a very basic knowledge of R is required.

Let's start loading essential packages for programming, customizing the outputs and defining a global ggplot theme.

## 1.2 Non-flexible models

These population dynamics models require at least two parameters, hence they are known as non-flexible, as opposed to the flexible ones for which there are at least one additional (third) parameter.

Following the convention proposed by (Madden et al. 2017a) in their book "The study of plant disease epidemics":

- time is represented by  $t$
- disease intensity by  $y$
- the rate of change in  $y$  between two time units is represented by  $\frac{dy}{dt}$

Now we can proceed and learn which non-flexible models exist and for which situation they are more appropriate.

### 1.2.1 Exponential

The differential equation for the exponential model is given by

$$\frac{dy}{dt} = r_E \cdot y,$$

where  $r_E$  is the apparent infection rate (subscript E for this model) (sensu Vanderplank) and  $y$  is the disease intensity. Biologically, this formulation suggests that diseased plants, or  $y$ , and  $r_E$  at each time contribute to disease increase. The value of  $\frac{dy}{dt}$  is minimal when  $y = 0$  and increases exponentially with the increase in  $y$ .

The integral for the exponential model is given by

$$y = y_0 e^{r_E t},$$

where  $y_0$  is and  $r$  are obtained via estimation. Let's simulate two curves by varying  $r$  while fixing  $y_0$  and varying the latter while fixing  $r_E$ . We produce the two plots in *ggplot* and add the predicted curve using the 'stat\_function'. But first, we need to define values for the two model

parameters. Further modifications to these values will be handled directly in the simulation (e.g. doubling infection rate, reducing initial inoculum by half, etc.).

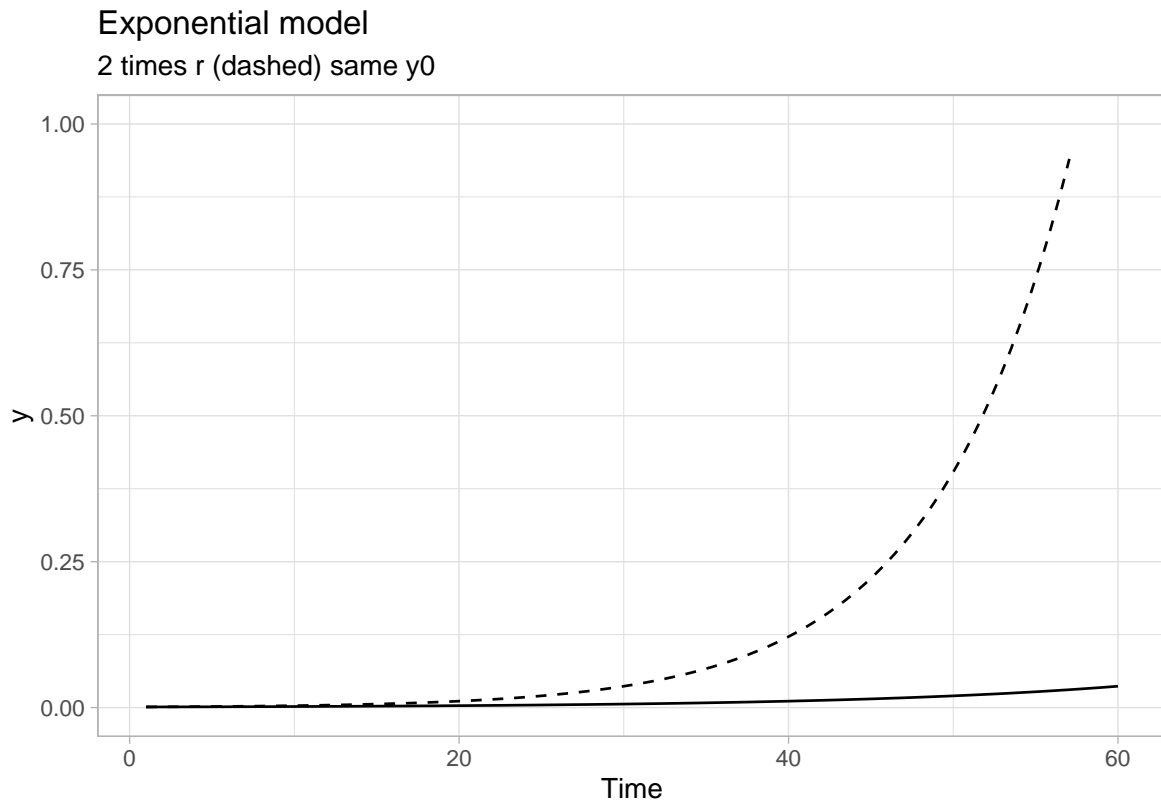
```
y0 <- 0.001
r <- 0.06
tmax <- 60 # maximum duration t of the epidemics
dat <- data.frame(t = seq(1:tmax), y = seq(0:1)) # define the axes
```

In the plot below, note that the infection rate in one curve was doubled ( $r = 0.12$ )

```
dat %>%
  ggplot(aes(t, y)) +
  stat_function(fun = function(t) y0 * exp(r * t), linetype = 1) +
  stat_function(fun = function(t) y0 * exp(r * 2 * t), linetype = 2) +
  ylim(0, 1) +
  labs(
    title = "Exponential model",
    subtitle = "2 times r (dashed) same y0",
    x = "Time"
  )
```

Warning: Removed 5 row(s) containing missing values (geom\_path).





Now the inoculum was increased five times while using the same doubled rate.

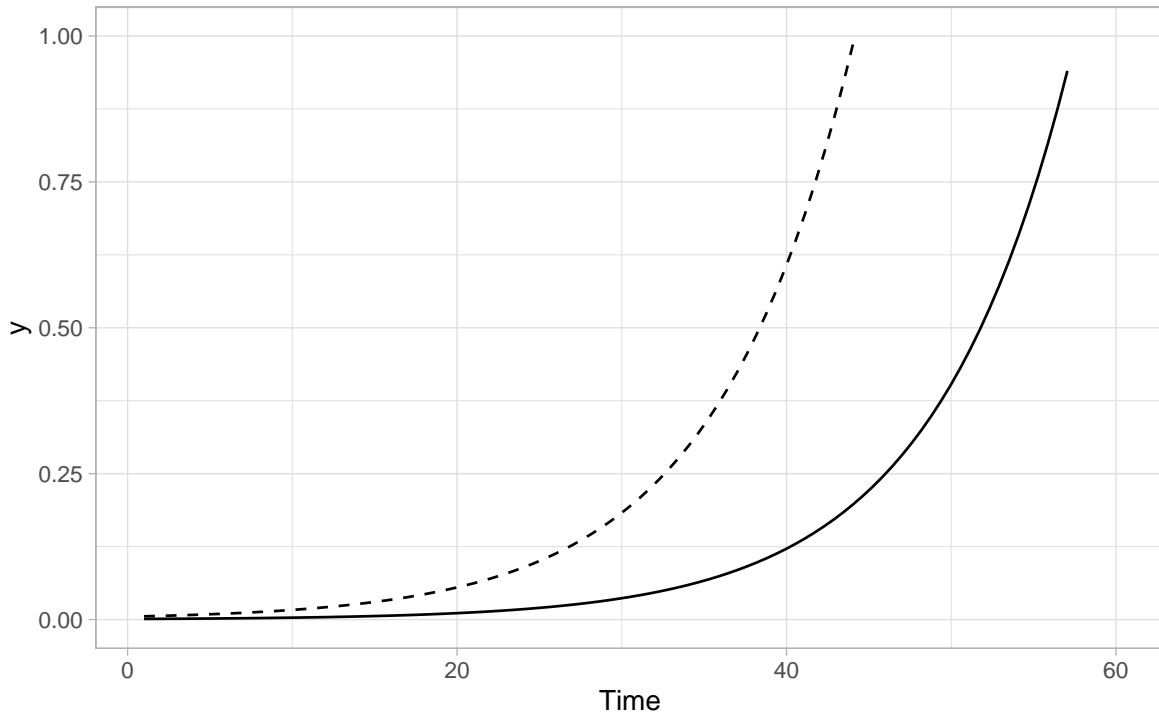
```
dat %>%  
  ggplot(aes(t, y)) +  
  stat_function(fun = function(t) y0 * exp(r * 2 * t), linetype = 1) +  
  stat_function(fun = function(t) y0 * 5 * exp(r * 2 * t), linetype = 2) +  
  ylim(0, 1) +  
  labs(title = "Exponential model", x = "Time",  
        subtitle = "5 times y0 (dashed) same r")
```

Warning: Removed 5 row(s) containing missing values (geom\_path).

Warning: Removed 27 row(s) containing missing values (geom\_path).

## Exponential model

5 times y0 (dashed) same r



### 1.2.2 Monomolecular

The differential of the monomolecular model is given by

$$\frac{dy}{dt} = r_M(1 - y)$$

where now the  $r_M$  is the rate parameter of the monomolecular model and  $(1 - y)$  is the proportion of non-infected (healthy) individuals or host tissue. Note that  $\frac{dy}{dt}$  is maximum when  $y = 0$  and decreases when  $y$  approaches 1. Its decline is due to decrease in the proportion of individuals or healthy sites with the increase in  $y$ . Any inoculum capable of infecting the host will more likely land on infected individuals or sites.

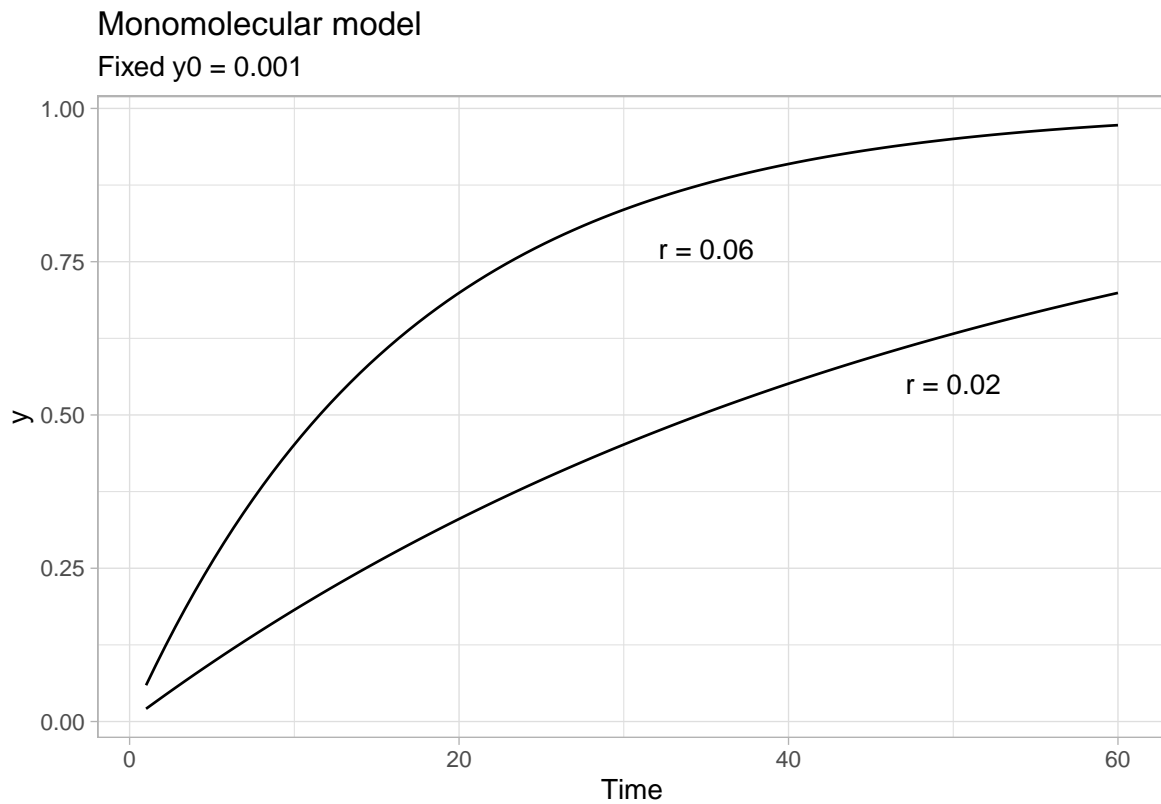
The integral of the monomolecular model is given by

$$\frac{dy}{dt} = 1 - (1 - y)e^{-r_M t}$$

This model commonly describes the temporal patterns of the monocyclic epidemics. In those, the inoculum produced during the course of the epidemics do not contribute new infections. Therefore, different from the exponential model, disease intensity  $y$  does not affect the epidemics and so the absolute rate is proportional to  $(1 - y)$ .

Let's simulate two monomolecular curve with different rate parameters where one is one third of the other.

```
dat %>%
  ggplot(aes(t, y)) +
  stat_function(fun = function(t) 1 - ((1 - y0) * exp(-r * t))) +
  stat_function(fun = function(t) 1 - ((1 - y0) * exp(-(r / 3) * t))) +
  labs(title = "Monomolecular model",
        subtitle = "Fixed y0 = 0.001", x = "Time")
  ) +
  annotate(geom = "text", x = 35, y = 0.77, label = "r = 0.06") +
  annotate(geom = "text", x = 50, y = 0.55, label = "r = 0.02")
```

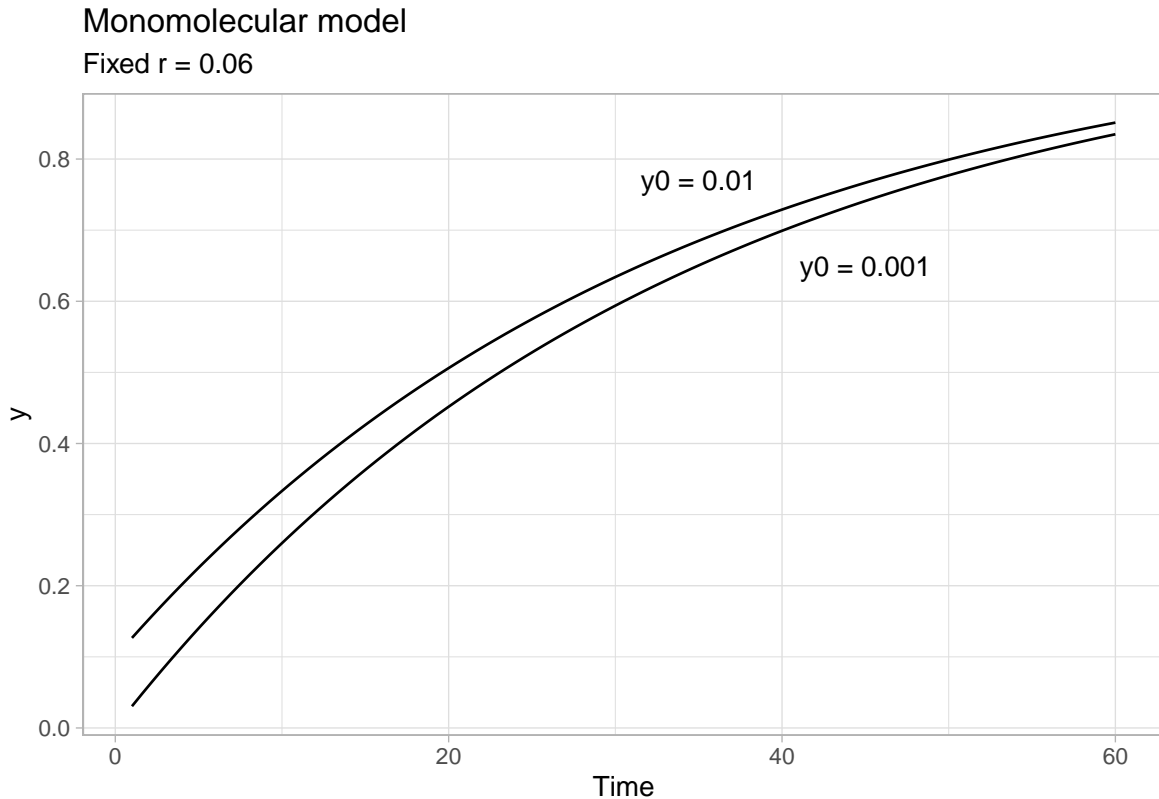


Now inoculum was increased 100 times with the reduced rate.

```

dat %>%
  ggplot(aes(t, y)) +
  stat_function(fun = function(t) 1 - ((1 - y0) * exp(-r / 2 * t))) +
  stat_function(fun = function(t) 1 - ((1 - (y0 * 100)) * exp(-r / 2 * t))) +
  labs(title = "Monomolecular model",
        subtitle = "Fixed r = 0.06", x = "Time") +
  annotate(geom = "text", x = 35, y = 0.77, label = "y0 = 0.01") +
  annotate(geom = "text", x = 45, y = 0.65, label = "y0 = 0.001")

```



### 1.2.3 Logistic

The logistic model is a more elaborated version of the two previous models as it incorporates the features of them both. Its differential is given by

$$\frac{dy}{dt} = r_L \cdot y \cdot (1 - y),$$

where  $r_L$  is the infection rate of the logistic model,  $y$  is the proportion of diseased individuals or host tissue and  $(1 - y)$  is the proportion of non-affected individuals or host area.

Biologically,  $y$  in its differential equation implies that  $\frac{dy}{dt}$  increases with the increase in  $y$  (as in the exponential) because more disease means more inoculum. However,  $(1 - y)$  leads to a decrease in  $\frac{dy}{dt}$  when  $y$  approaches the maximum  $y = 1$ , because the proportion of healthy individuals or host area decreases (as in the monomolecular). Therefore,  $\frac{dy}{dt}$  is minimal at the onset of the epidemics, reaches a maximum when  $y/2$  and declines until  $y = 1$ .

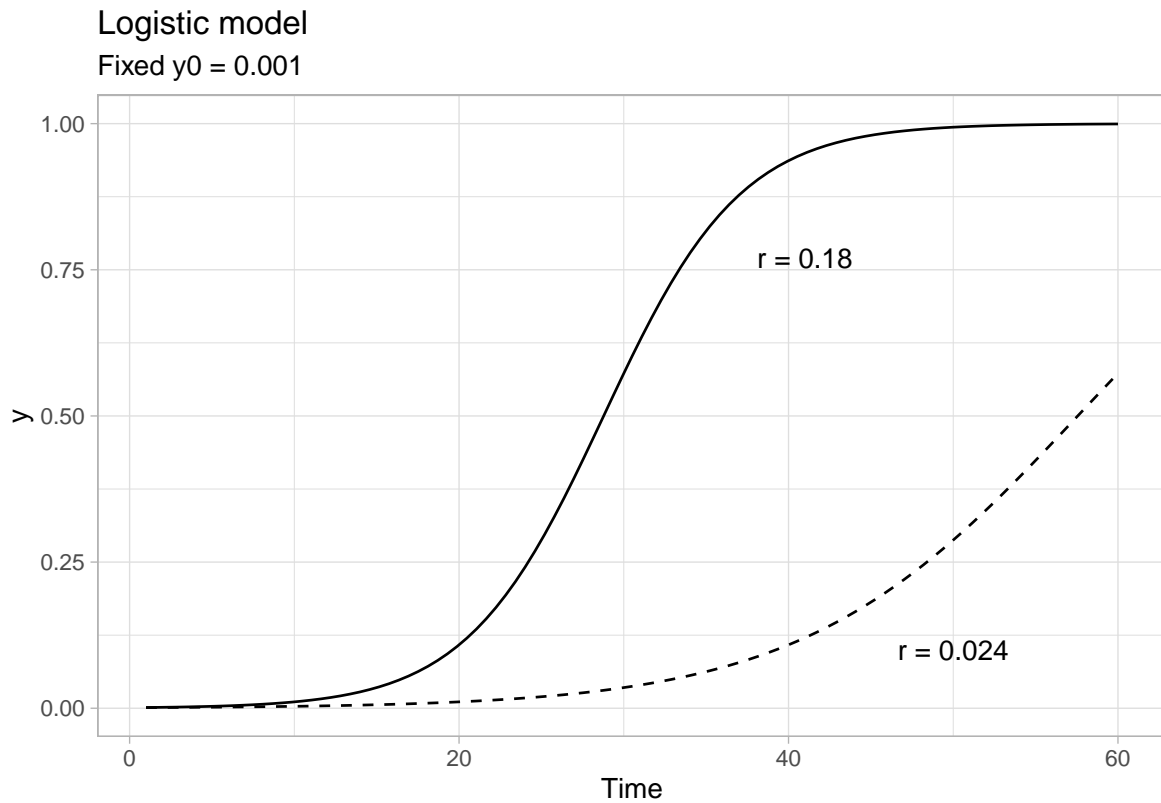
The integral is given by

$$y = \frac{1}{1 + (1 - y_0) \cdot e^{-r \cdot t}},$$

where  $r_L$  is the apparent infection rate of the logistic model and  $y_0$  is the disease intensity at  $t = 0$ . This model provides a good fit to polycyclic epidemics.

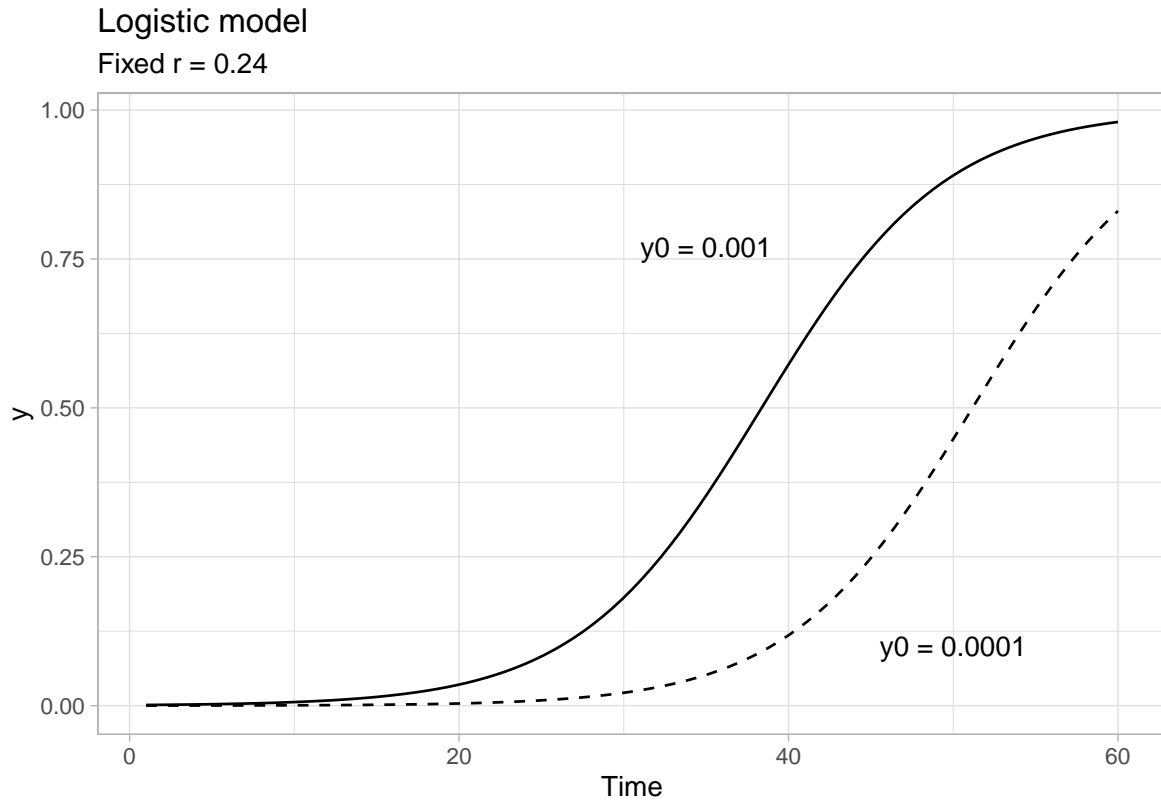
Let's check two curves where in one the infection rate is double while keeping the same initial inoculum.

```
dat %>%
  ggplot(aes(t, y)) +
  stat_function(
    linetype = 2,
    fun = function(t) 1 / (1 + ((1 - y0) / y0) * exp(-r * 2 * t))
  ) +
  stat_function(fun = function(t) 1 / (1 + ((1 - y0) / y0) * exp(-r * 4 * t))) +
  labs(title = "Logistic model", subtitle = "Fixed y0 = 0.001", x = "Time") +
  annotate(geom = "text", x = 41, y = 0.77, label = "r = 0.18") +
  annotate(geom = "text", x = 50, y = 0.10, label = "r = 0.024")
```



Now the inoculum is reduced 10 times for a same infection rate.

```
dat %>%
  ggplot(aes(t, y)) +
  stat_function(
    linetype = 2,
    fun = function(t) 1 / (1 + ((1 - (y0 / 10)) / (y0 / 10)) * exp(-r * 3 * t))
  ) +
  stat_function(fun = function(t) 1 / (1 + ((1 - y0) / y0) * exp(-r * 3 * t))) +
  labs(title = "Logistic model", subtitle = "Fixed r = 0.24", x = "Time") +
  annotate(geom = "text", x = 35, y = 0.77, label = "y0 = 0.001") +
  annotate(geom = "text", x = 50, y = 0.10, label = "y0 = 0.0001")
```



#### 1.2.4 Gompertz

The Gompertz model is similar to the logistic and also provides a very good fit to several polycyclic diseases. The differential equation is given by

$$\frac{dy}{dt} = r_G \cdot [\ln(1) - \ln(y)]$$

Differently from the logistic, the variable representing the non-infected individuals or host area is  $-\ln(y)$ . The integral equation is given by

$$y = e^{(\ln(y_0)) \cdot e^{-r_G \cdot t}},$$

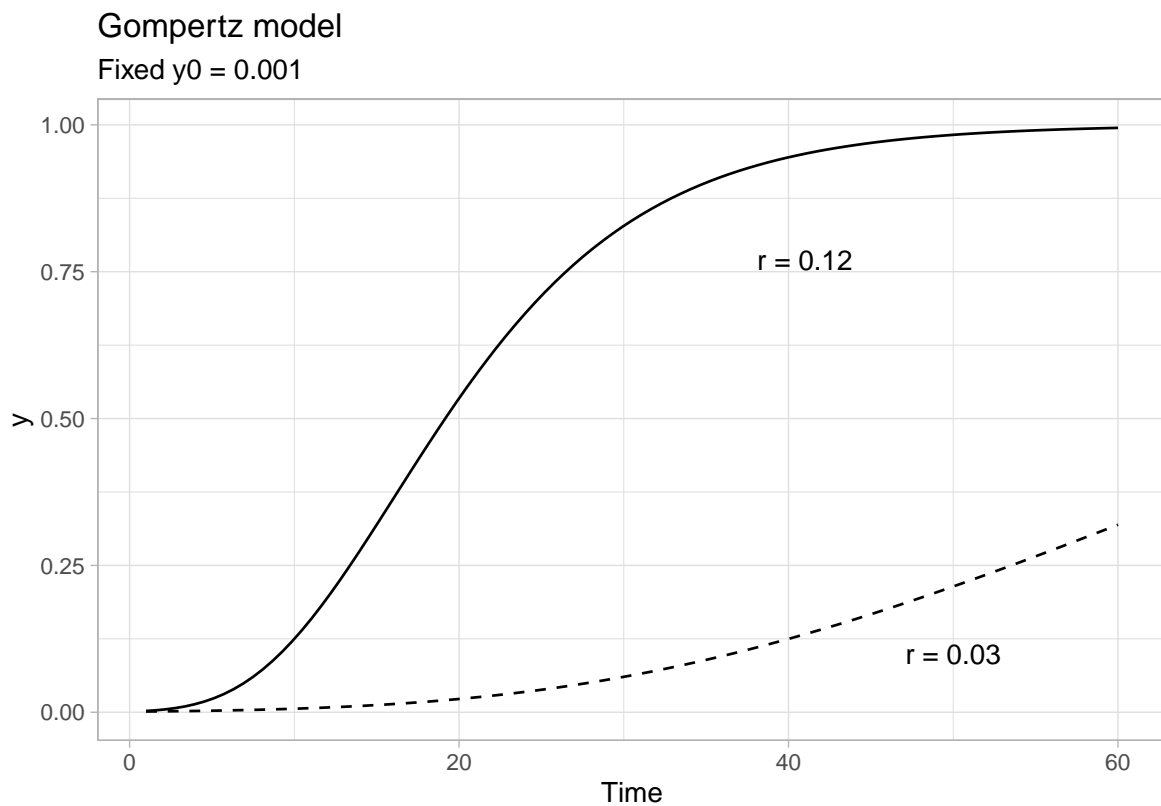
where  $r_G$  is the apparent infection rate for the Gompertz models and  $y_0$  is the disease intensity at  $t = 0$ .

Let's check curves for two rates.

```

dat %>%
  ggplot(aes(t, y)) +
  stat_function(
    linetype = 2,
    fun = function(t) exp(log(y0) * exp(-r/2 * t))
  ) +
  stat_function(fun = function(t) exp(log(y0) * exp(-r*2 * t))) +
  labs(title = "Gompertz model", subtitle = "Fixed y0 = 0.001", x = "Time") +
  annotate(geom = "text", x = 41, y = 0.77, label = "r = 0.12") +
  annotate(geom = "text", x = 50, y = 0.10, label = "r = 0.03")

```



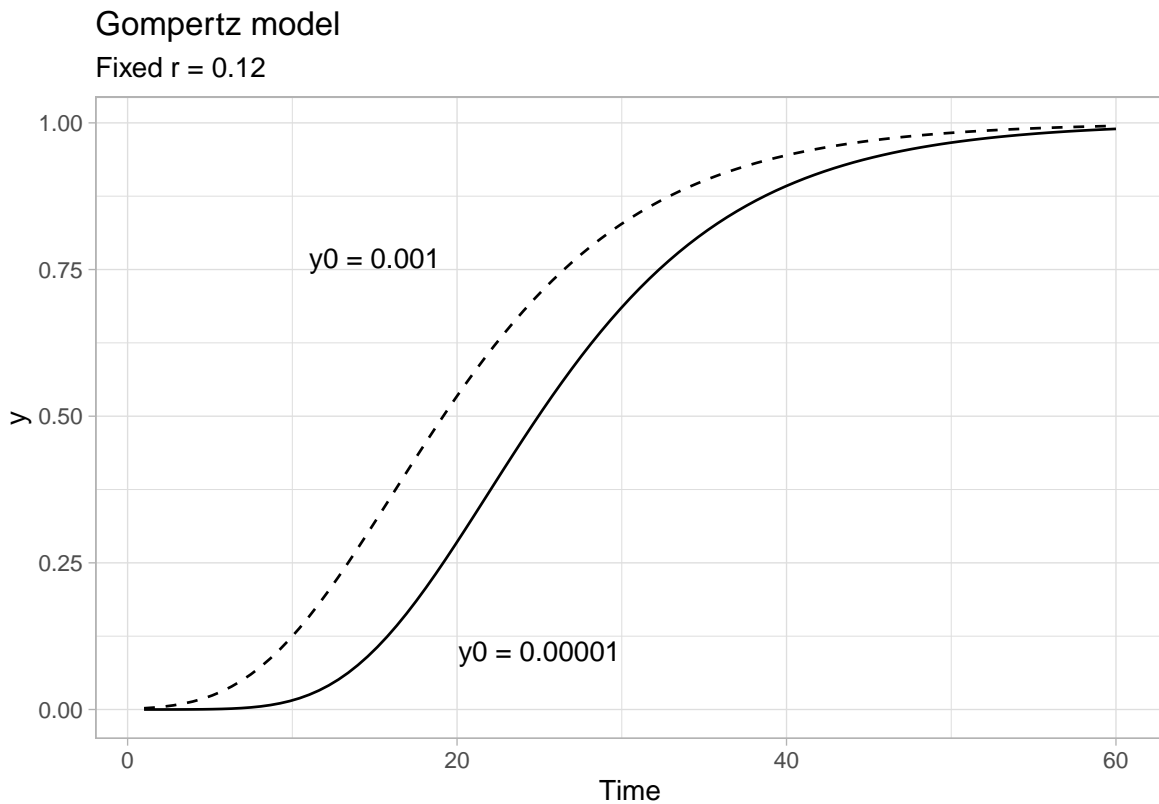
And those when inoculum was reduced thousand times.



```

dat %>%
  ggplot(aes(t, y)) +
  stat_function(
    linetype = 2,
    fun = function(t) exp(log(y0) * exp(-r*2 * t))
  ) +
  stat_function(fun = function(t) exp(log(y0/1000) * exp(-r*2 * t))) +
  labs(title = "Gompertz model", subtitle = "Fixed r = 0.12", x = "Time") +
  annotate(geom = "text", x = 15, y = 0.77, label = "y0 = 0.001") +
  annotate(geom = "text", x = 25, y = 0.10, label = "y0 = 0.00001")

```



### 1.3 Model fitting

In this tutorial you will learn how to fit models to multiple actual disease progress curves (DPCs) data obtained from the literature. I will demonstrate how to fit and select the models

using a new R package called *epifitter*. A few user friendly functions will help us decide which model to choose to obtain the parameters of interest and further compare the epidemics.

To illustrate, I will use two datasets available from Chapter 3 from the book, *Study of Plant Disease Epidemics* (Madden et al. 2017a). In the book, SAS codes are presented to perform a few analysis. We then provide an alternative code for performing similar analysis, although not perfectly reproducing the results from the book.

### 1.3.1 Non-replicated

Here we will compare three DPCs of the incidence of tobacco etch, a virus disease, in peppers. Evaluations of incidence were evaluated at a 7-day interval up to 49 days. The data are available in chapter 4 (page 93). Let's input the data manually and create a data frame. First column is the assessment time and the other columns correspond to the treatments, called groups in the book, from 1 to 3.

#### 1.3.1.1 Initial setup

Load essential packages and set parameters recursively.

```
library(tidyverse)
library(knitr)
library(patchwork)
library(ggthemes)
theme_set(theme_few())
knitr::opts_chunk$set(echo=TRUE, warning=FALSE, message=FALSE)
options(digits = 3)
```

### 1.3.1.2 Entering data

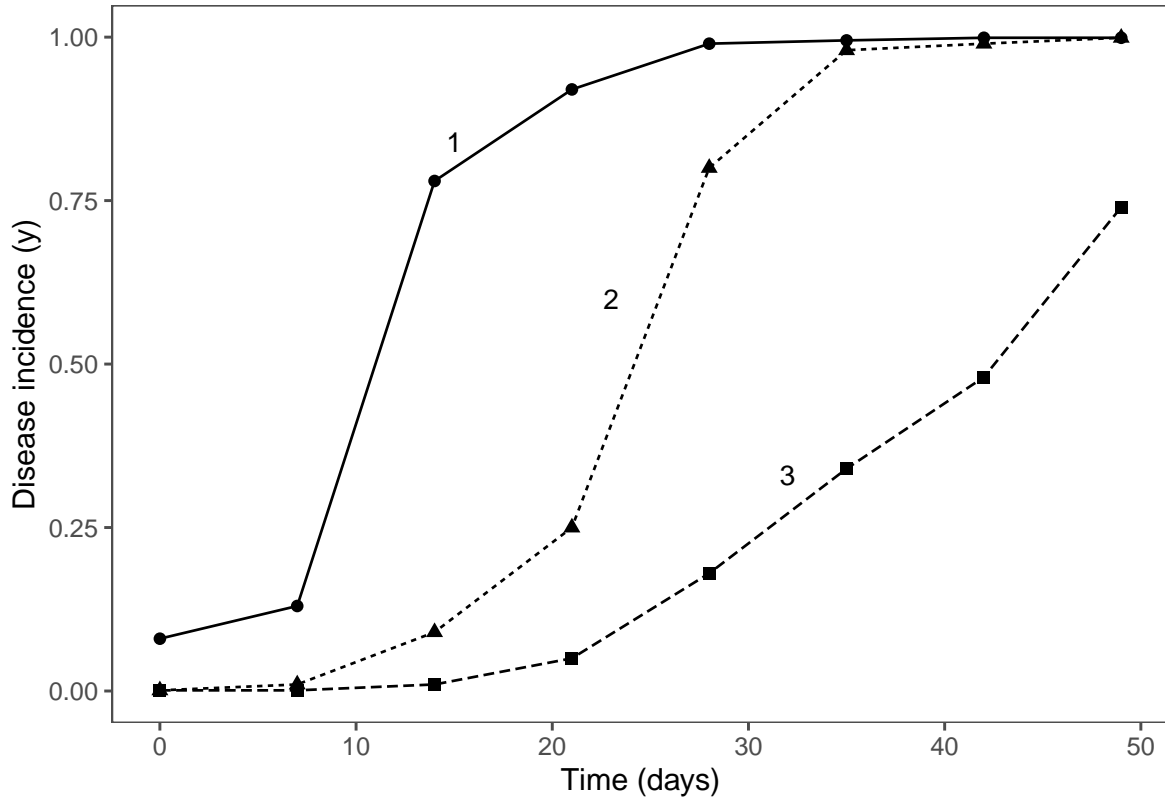
```
pepper <-  
  tibble::tribble(  
    ~t, ~`1`, ~`2`, ~`3`,  
    0, 0.08, 0.001, 0.001,  
    7, 0.13, 0.01, 0.001,  
    14, 0.78, 0.09, 0.01,  
    21, 0.92, 0.25, 0.05,  
    28, 0.99, 0.8, 0.18,  
    35, 0.995, 0.98, 0.34,  
    42, 0.999, 0.99, 0.48,  
    49, 0.999, 0.999, 0.74  
  )
```

### 1.3.1.3 Visualize the DPCs

Before proceeding with model selection and fitting, let's visualize the three epidemics. The code below reproduces quite exactly the top plot of Fig. 4.15 ((Madden et al. 2017a) page 94). The appraisal of the curves might give us a hint on which models are the best candidates.

Because the data was entered in the wide format (each DPCs in a different columns) we need to reshape it to the tidyverse-suitable format, which is the long format. The `pivot_longer` function will do the job of reshaping from wide to long format so we can finally use the `ggplot` function to produce the plot.

```
pepper %>%  
  pivot_longer(2:4, names_to = "treat", values_to = "inc") %>%  
  ggplot (aes(t, inc,  
              linetype = treat,  
              shape = treat,  
              group = treat))+  
  geom_point(size =2)+  
  geom_line()+  
  annotate(geom = "text", x = 15, y = 0.84, label = "1")+  
  annotate(geom = "text", x = 23, y = 0.6, label = "2")+  
  annotate(geom = "text", x = 32, y = 0.33, label = "3")+  
  labs(y = "Disease incidence (y)",  
       x = "Time (days)")+  
  theme(legend.position = "none")
```



Most of the three curves show a sigmoid shape with the exception of group 3 that resembles an exponential growth, not reaching the maximum value, and thus suggesting an incomplete epidemic. We can easily eliminate the monomolecular and exponential models and decide on the other two non-flexible models: logistic or Gompertz. To do that, let's proceed to model fitting and evaluate the statistics for supporting a final decision. There are two modeling approaches for model fitting in *epifitter*: the **linear** or **nonlinear** parameter-estimation methods.

#### 1.3.1.4 Fitting: single epidemics

Among the several options offered by *epifitter* we start with the simplest one, which is fit a model to a single epidemics using the linear regression approach. For such, the `fit_lin()` requires two arguments: time (`time`) and disease intensity (`y`) each one as a vector stored or not in a dataframe.

Since we have three epidemics, `fit_lin()` will be use three times. The function produces a list object with six elements. Let's first look at the `Stats` dataframe of each of the three lists named `epi1` to `epi3`.

```
library(epifitter)
epi1 <- fit_lin(time = pepper$t,
                y = pepper$`1` )
epi1$Stats
```

	CCC	r_squared	RSE
Gompertz	0.985	0.970	0.591
Monomolecular	0.984	0.968	0.543
Logistic	0.978	0.957	0.824
Exponential	0.784	0.645	0.670

```
epi2 <- fit_lin(time = pepper$t,
                y = pepper$`2` )
epi2$Stats
```

	CCC	r_squared	RSE
Logistic	0.996	0.992	0.452
Gompertz	0.971	0.943	0.841
Monomolecular	0.925	0.860	1.068
Exponential	0.897	0.813	1.202

```
epi3 <- fit_lin(time = pepper$t,
                y = pepper$`3` )
epi3$Stats
```

	CCC	r_squared	RSE
Logistic	0.983	0.967	0.605
Gompertz	0.983	0.966	0.226
Exponential	0.964	0.930	0.771
Monomolecular	0.859	0.753	0.253

The statistics of the model fit confirms our initial guess that the predictions by the logistic or the Gompertz are closer to the observations than predictions by the other models. There is no much difference between them based on these statistics. However, to pick one of the models, it is important to inspect the curves with the observed and predicted values to check which model is best for all curves.

### 1.3.1.5 Fitting: multiple epidemics

Before looking at the prediction, let's use another handy function that allows us to simultaneously fit the models to multiple DPC data. Different from `fit_lin()`, `fit_multi()` requires the data to be structured in the long format where there is a column specifying each of the epidemics.

Let's then create a new data set called `pepper2` using the data transposing functions of the *tidyr* package.

```
pepper2 <- pepper %>%  
  pivot_longer(2:4, names_to = "treat", values_to = "inc")
```

Now we fit the models to all DPCs. Note that the name of the variable indicating the DPC code needs to be informed in `strata_cols` argument.

```
epi_all <- fit_multi(  
  time_col = "t",  
  intensity_col = "inc",  
  data = pepper2,  
  strata_cols = "treat",  
  nlin = FALSE  
)
```

Now let's select the statistics of model fitting. Again, *Epifitter* ranks the models based on the CCC (the higher the better) but it is important to check the RSE as well - the lower the better. In fact, the RSE is more important when the goal is prediction.

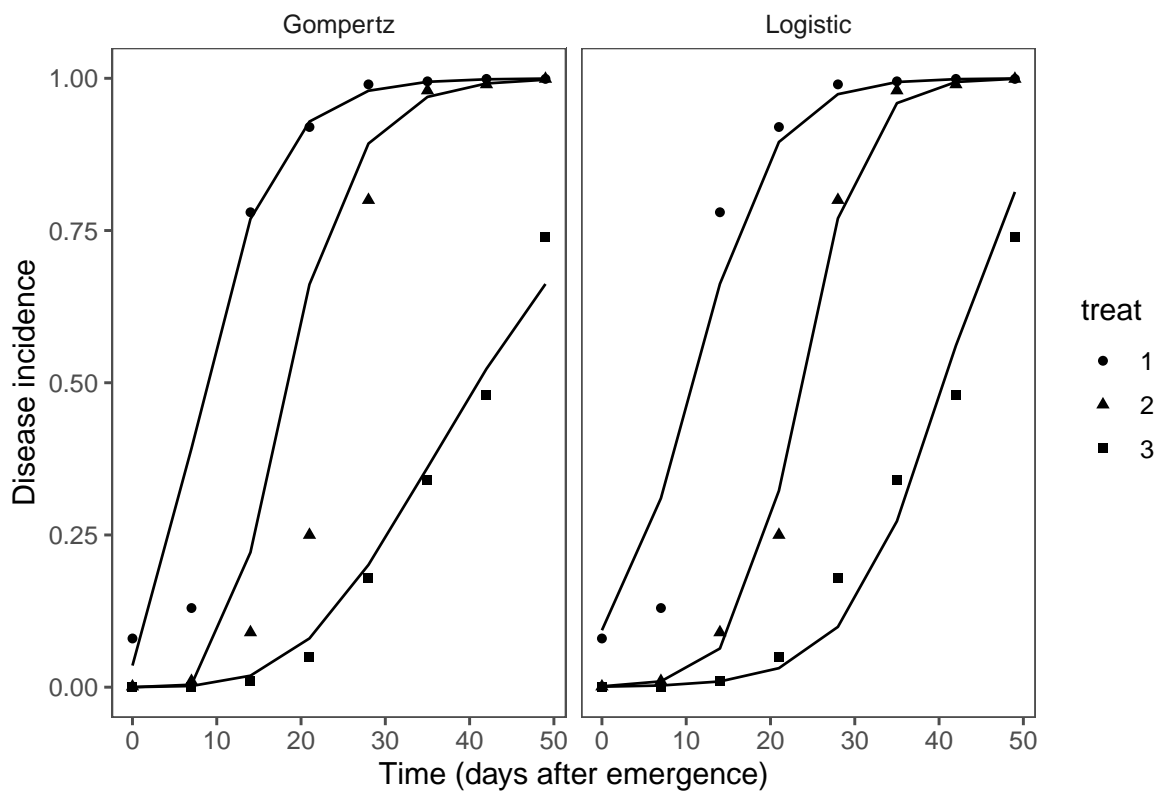
```
epi_all$Parameters %>%  
  select(treat, model, best_model, RSE, CCC)
```

	treat	model	best_model	RSE	CCC
1	1	Gompertz	1	0.591	0.985
2	1	Monomolecular	2	0.543	0.984
3	1	Logistic	3	0.824	0.978
4	1	Exponential	4	0.671	0.784
5	2	Logistic	1	0.452	0.996
6	2	Gompertz	2	0.841	0.971
7	2	Monomolecular	3	1.068	0.925
8	2	Exponential	4	1.202	0.897
9	3	Logistic	1	0.605	0.983
10	3	Gompertz	2	0.226	0.982

11	3	Exponential	3	0.771	0.964
12	3	Monomolecular	4	0.253	0.859

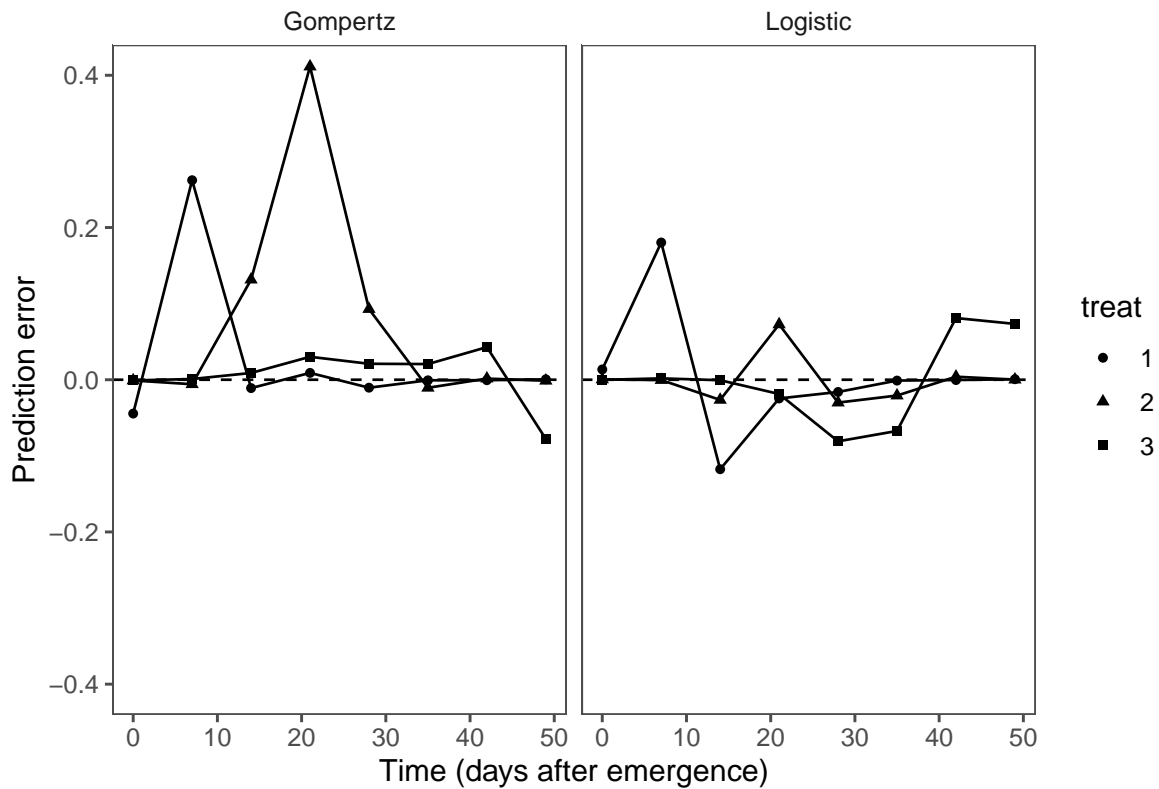
To be more certain about our decision, let's advance to the final step which is to produce the plots with the observed and predicted values for each assessment time by calling the `Data` dataframe of the `'epi_all'` list.

```
epi_all$Data %>%
  filter(model %in% c("Gompertz", "Logistic")) %>%
  ggplot(aes(time, predicted, shape = treat)) +
  geom_point(aes(time, y)) +
  geom_line() +
  facet_wrap(~ model) +
  coord_cartesian(ylim = c(0, 1)) + # set the max to 0.6
  labs(
    y = "Disease incidence",
    x = "Time (days after emergence)"
  )
)
```



Overall, the logistic model seems a better fit for all the curves. Let's produce a plot with the prediction error versus time.

```
epi_all$Data %>%
  filter(model %in% c("Gompertz", "Logistic")) %>%
  ggplot(aes(time, predicted -y, shape = treat)) +
  geom_point() +
  geom_line() +
  geom_hline(yintercept = 0, linetype = 2) +
  facet_wrap(~ model) +
  coord_cartesian(ylim = c(-0.4, 0.4)) + # set the max to 0.6
  labs(
    y = "Prediction error",
    x = "Time (days after emergence)"
  )
)
```



The plots above confirms the logistic model as good fit overall because the errors for all epidemics combined are more scattered around the non-error line.



```

epi_all$Parameters %>%
  filter(model == "Logistic") %>%
  select(treat, y0, y0_ci_lwr, y0_ci_upr, r, r_ci_lwr, r_ci_upr
)

```

	treat	y0	y0_ci_lwr	y0_ci_upr	r	r_ci_lwr	r_ci_upr
1	1	0.093504	0.027321	0.27473	0.210	0.166	0.255
2	2	0.001373	0.000672	0.00280	0.278	0.254	0.303
3	3	0.000813	0.000313	0.00211	0.175	0.143	0.208

We can produce a plot for visual inference on the differences in the parameters.

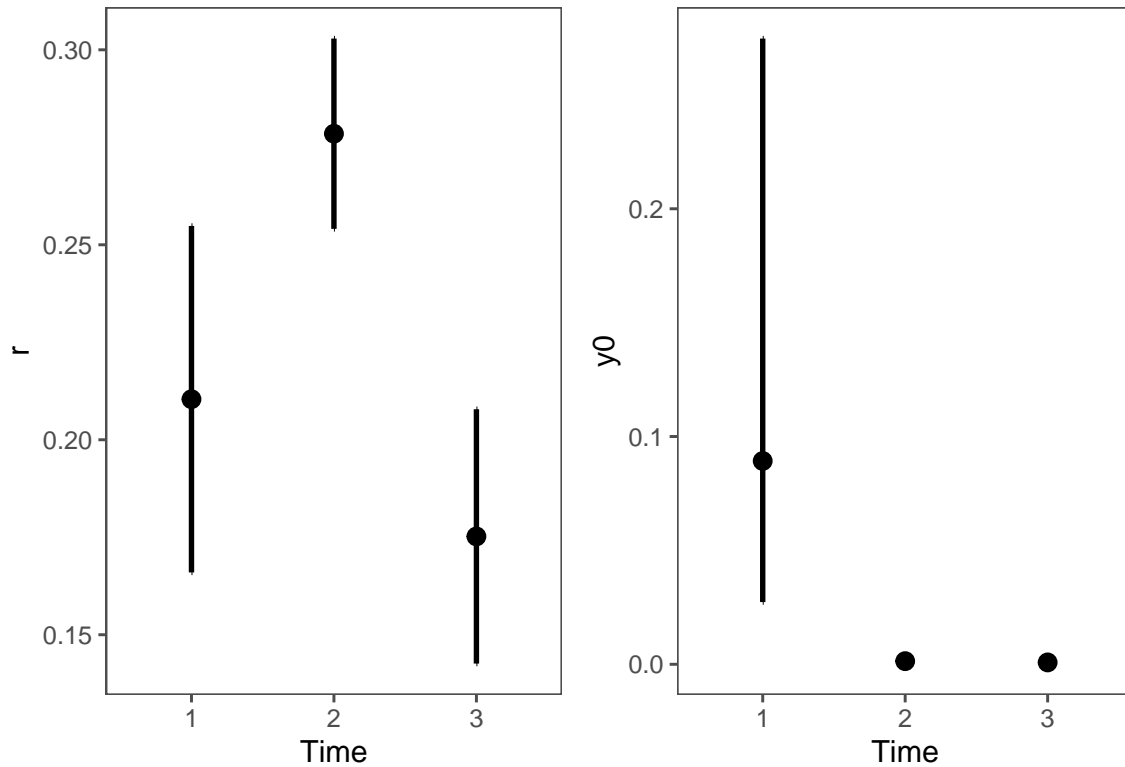
```

p1 <- epi_all$Parameters %>%
  filter(model == "Logistic") %>%
  ggplot(aes(treat, r)) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = r_ci_lwr, ymax = r_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "r"
  )

p2 <- epi_all$Parameters %>%
  filter(model == "Logistic") %>%
  ggplot(aes(treat, 1 - exp(-y0))) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = y0_ci_lwr, ymax = y0_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "y0"
  )

p1 | p2

```



### 1.3.2 Designed experiments

In this next section, we will work with disease data collected over time in the same plot unit (also called repeated measures) from a designed experiment for evaluating and comparing treatment effects.

Again, we will use a dataset of progress curves shown in page 98 (Madden et al. 2017a). The curves represent the incidence of soybean plants symptomatic for bud blight caused by tobacco streak virus. Four treatments (different planting dates) were evaluated in randomized complete block design with four replicates. There are four assessment in time for each curve. The data was stored as a csv file and will be loaded using `read_csv()` function and stored as dataframe called `budblight`.

#### 1.3.2.1 Loading data

```
budblight <- read_csv("data/bud-blight-soybean.csv")
```

Let's have a look at the first six rows of the dataset and check the data type for each column. There is an additional column representing the replicates, called block.

```
head(budblight)
```

```
# A tibble: 6 x 4
  treat time block    y
  <chr> <dbl> <dbl> <dbl>
1 PD1    30     1  0.1
2 PD1    30     2  0.3
3 PD1    30     3  0.1
4 PD1    30     4  0.1
5 PD1    40     1  0.3
6 PD1    40     2  0.38
```

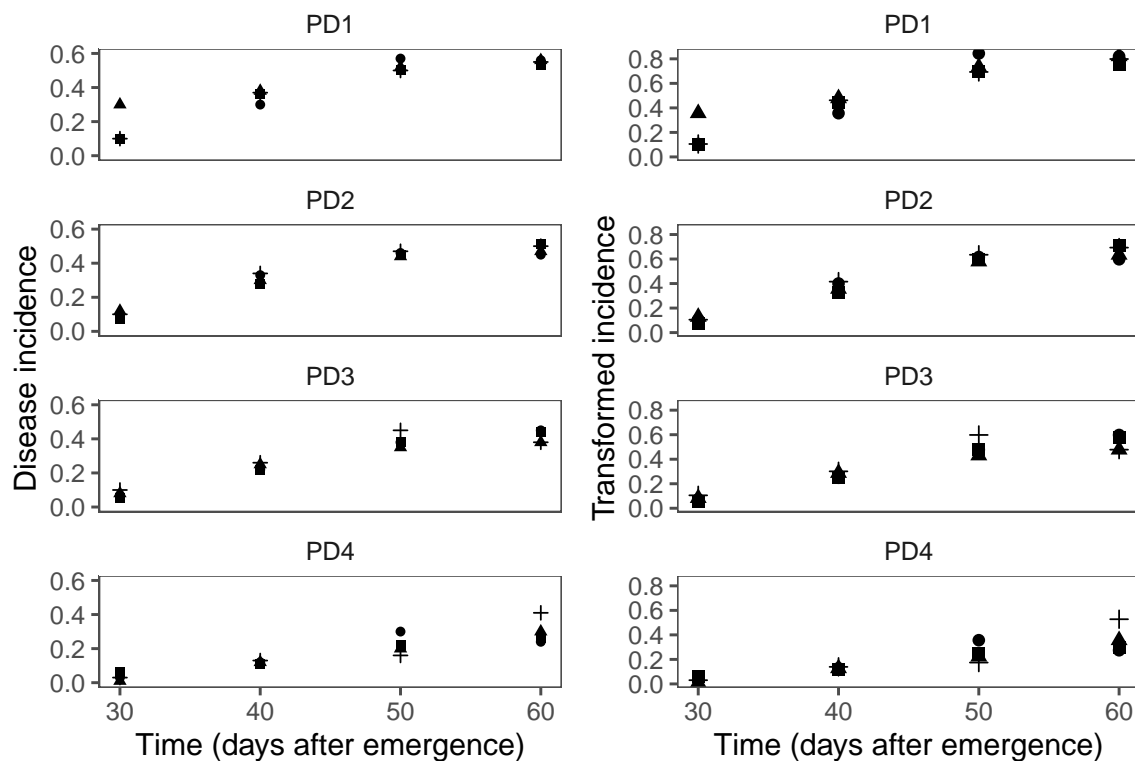
### 1.3.2.2 Visualizing the DPCs

Let's have a look at the curves and produce a combo plot figure similar to Fig. 4.17 of the book, but without the line of the predicted values.

```
p3 <- budblight %>%
  ggplot(aes(
    time, y,
    group = block,
    shape = factor(block)
  )) +
  geom_point(size = 1.5) +
  ylim(0, 0.6) +
  theme(legend.position = "none")+
  facet_wrap(~treat, ncol = 1)+
  labs(y = "Disease incidence",
       x = "Time (days after emergence)")
```

```
p4 <- buddblight %>%
  ggplot(aes(
    time, log(1 / (1 - y)),
    group = block,
    shape = factor(block)
  )) +
  geom_point(size = 2) +
  facet_wrap(~treat, ncol = 1) +
  theme(legend.position = "none")+
  labs(y = "Transformed incidence", x = "Time (days after emergence)")
```

p3 | p4



### 1.3.2.3 Model fitting

Remember that the first step in model selection is the visual appraisal of the curve data linearized with the model transformation. In the case the curves represent complete epidemics

(close to 100%) appraisal of the absolute rate (difference in  $y$  between two times) over time is also helpful.

For the treatments above, it looks like the curves are typical of a monocyclic disease (the case of soybean bud blight), for which the monomolecular is usually a good fit, but other models are also possible as well. For this exercise, we will use both the linear and the nonlinear estimation method.

### 1.3.2.3.1 Linear regression

For convenience, we use the `fit_multi()` to handle multiple epidemics. The function returns a list object where a series of statistics are provided to aid in model selection and parameter estimation. We need to provide the names of columns (arguments): assessment time (`time_col`), disease incidence (`intensity_col`), and treatment (`strata_cols`).

```
lin1 <- fit_multi(
  time_col = "time",
  intensity_col = "y",
  data = budblight,
  strata_cols = "treat",
  nlin = FALSE
)
```

Let's look at how well the four models fitted the data. Epifitter suggests the best fitted model (1 to 4, where 1 is best) for each treatment. Let's have a look at the statistics of model fitting.

```
lin1$Parameters %>%
  select(treat, best_model, model, CCC, RSE)
```

	treat	best_model	model	CCC	RSE
1	PD1	1	Monomolecular	0.935	0.0981
2	PD1	2	Gompertz	0.904	0.2223
3	PD1	3	Logistic	0.871	0.4475
4	PD1	4	Exponential	0.828	0.3612
5	PD2	1	Monomolecular	0.955	0.0700
6	PD2	2	Gompertz	0.931	0.1794
7	PD2	3	Logistic	0.906	0.3877
8	PD2	4	Exponential	0.880	0.3268
9	PD3	1	Monomolecular	0.939	0.0683
10	PD3	2	Gompertz	0.929	0.1716
11	PD3	3	Logistic	0.909	0.3905

12	PD3	4	Exponential	0.890	0.3388
13	PD4	1	Gompertz	0.923	0.1747
14	PD4	2	Monomolecular	0.895	0.0649
15	PD4	3	Logistic	0.891	0.5241
16	PD4	4	Exponential	0.874	0.4977

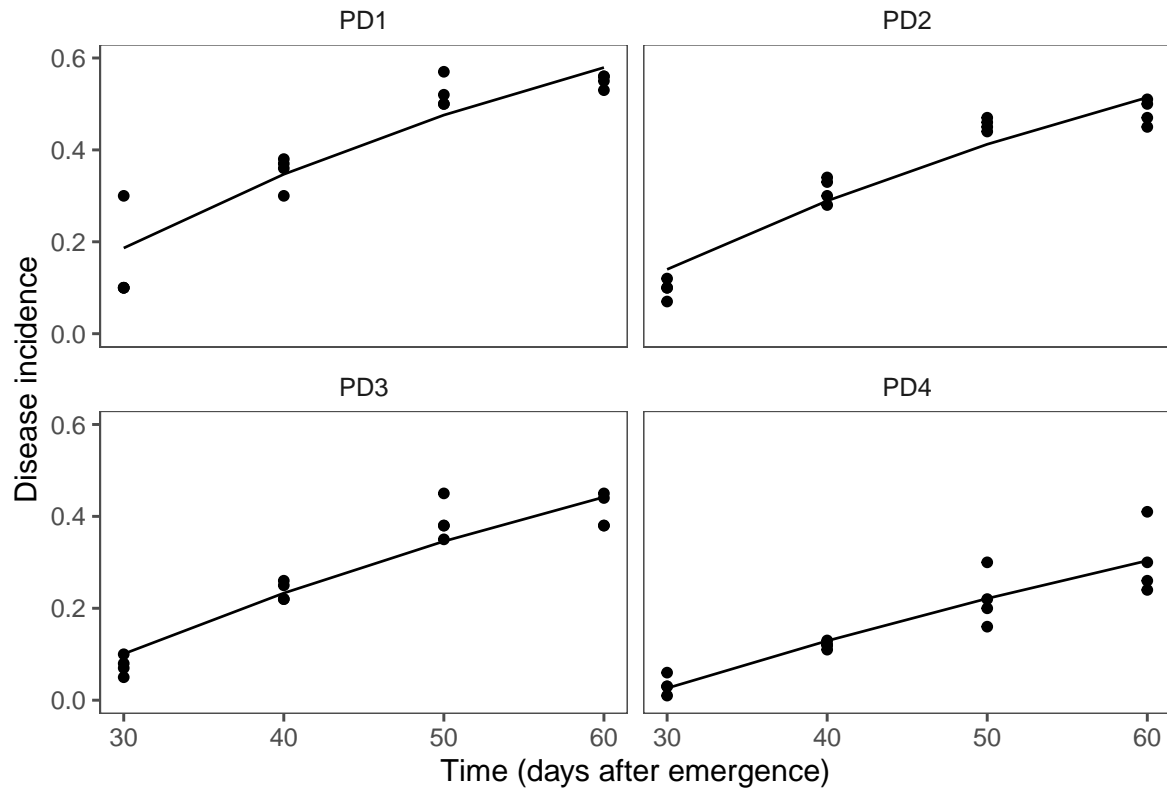
And now we extract values for each parameter estimated from the fit of the monomolecular model.

```
lin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  select(treat, y0, r)
```

	treat	y0	r
1	PD1	-0.573	0.0220
2	PD2	-0.522	0.0190
3	PD3	-0.449	0.0159
4	PD4	-0.362	0.0112

Now we visualize the fit of the monomolecular model (using `filter` function - see below) to the data together with the observed data and then reproduce the right plots in Fig. 4.17 from the book.

```
lin1$Data %>%
  filter(model == "Monomolecular") %>%
  ggplot(aes(time, predicted)) +
  geom_point(aes(time, y)) +
  geom_line(size = 0.5) +
  facet_wrap(~treat) +
  coord_cartesian(ylim = c(0, 0.6)) + # set the max to 0.6
  labs(
    y = "Disease incidence",
    x = "Time (days after emergence)"
  )
```



Now we can plot the means and respective 95% confidence interval of the apparent infection rate ( $r$ ) and initial inoculum ( $y_0$ ) for visual inference.

```

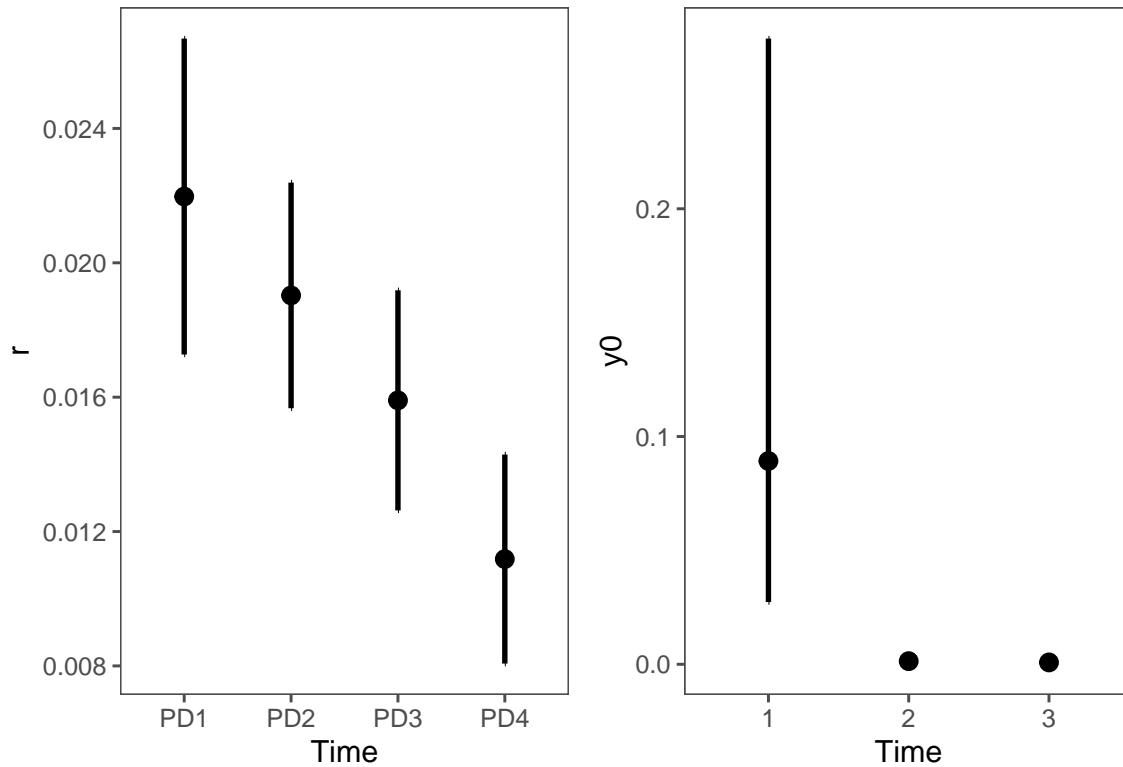
p5 <- lin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  ggplot(aes(treat, r)) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = r_ci_lwr, ymax = r_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "r"
  )

p6 <- lin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  ggplot(aes(treat, 1 - exp(-y0))) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = y0_ci_lwr, ymax = y0_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "y0"
  )

p5 | p2

```





#### 1.3.2.3.2 Non-linear regression

To estimate the parameters using the non-linear approach, we repeat the same arguments in the `fit_multi` function, but include an additional argument `nlin` set to `TRUE`.

```
nlin1 <- fit_multi(
  time_col = "time",
  intensity_col = "y",
  data = budblight,
  strata_cols = "treat",
  nlin = TRUE
)
```

Let's check statistics of model fit.

```
nlin1$Parameters %>%
  select(treat, model, CCC, RSE, best_model)
```

	treat	model	CCC	RSE	best_model
1	PD1	Monomolecular	0.938	0.0613	1
2	PD1	Gompertz	0.917	0.0699	2
3	PD1	Logistic	0.896	0.0770	3
4	PD1	Exponential	0.854	0.0880	4
5	PD2	Monomolecular	0.967	0.0421	1
6	PD2	Gompertz	0.935	0.0573	2
7	PD2	Logistic	0.908	0.0666	3
8	PD2	Exponential	0.870	0.0767	4
9	PD3	Monomolecular	0.957	0.0427	1
10	PD3	Gompertz	0.926	0.0544	2
11	PD3	Logistic	0.900	0.0620	3
12	PD3	Exponential	0.870	0.0689	4
13	PD4	Monomolecular	0.918	0.0460	1
14	PD4	Gompertz	0.909	0.0479	2
15	PD4	Logistic	0.894	0.0508	3
16	PD4	Exponential	0.884	0.0527	4

And now we obtain the two parameters of interest. Note that the values are not the same as those estimated using linear regression, but they are similar and highly correlated.

```
nlin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  select(treat, y0, r)
```

	treat	y0	r
1	PD1	-0.707	0.0238
2	PD2	-0.634	0.0206
3	PD3	-0.505	0.0167
4	PD4	-0.350	0.0109

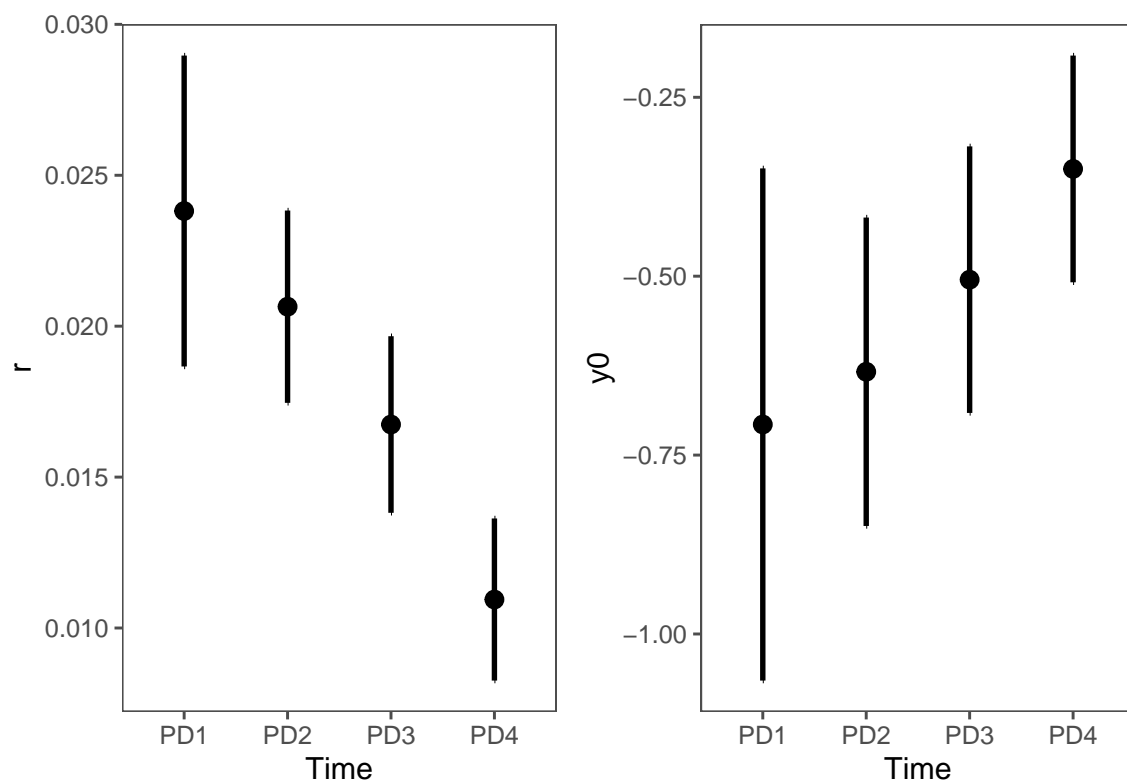
```

p7 <- nlin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  ggplot(aes(treat, r)) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = r_ci_lwr, ymax = r_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "r"
  )

p8 <- nlin1$Parameters %>%
  filter(model == "Monomolecular") %>%
  ggplot(aes(treat, y0)) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = y0_ci_lwr, ymax = y0_ci_upr),
    width = 0,
    size = 1
  ) +
  labs(
    x = "Time",
    y = "y0"
  )

p7 | p8

```



## 2 Spatial gradients

### 2.1 Models

When modeling disease gradients, the distance is represented by  $x$ , a continuous variable which can be expressed by various units (cm, m, km, etc). The gradient models, similar to the population dynamics models (disease progress) are of the **deterministic** type. The difference is that, for disease progress curves, disease intensity tends to increase with increasing time, while in disease gradients the disease intensity tends to decrease with increasing distance from the source of inoculum. Two models are most commonly fitted to data on disease gradients. More details about these models can be obtained it [this tutorial](#).

#### 2.1.1 Exponential model

The exponential model is also known as Kiyosawa & Shiyomi model. The differential of the exponential model is given by

$$\frac{dy}{dx} = -b_E \cdot y ,$$

where  $b_E$  is the exponential form of the rate of decline and  $y$  is the disease intensity. This model suggests that  $y$  (any disease intensity) is greater close to the source of inoculum, or at the distance zero. The integral form of the model is given by

$$y = a \cdot e^{-b \cdot x} ,$$

where  $a$  is the disease intensity at the distance zero and  $b$  is the rate of decline, in this case negative because disease intensity decreases with the increase of the distance from inoculum source. Let's make a plot for two disease gradients of varying parameters for this model.

First we need to load essential packages for programming, customizing the outputs and defining a global ggplot theme.

```
library(tidyverse)
library(ggthemes)
library(patchwork)
theme_set(theme_light())
knitr::opts_chunk$set(echo = TRUE, warning = FALSE, message = FALSE)
```

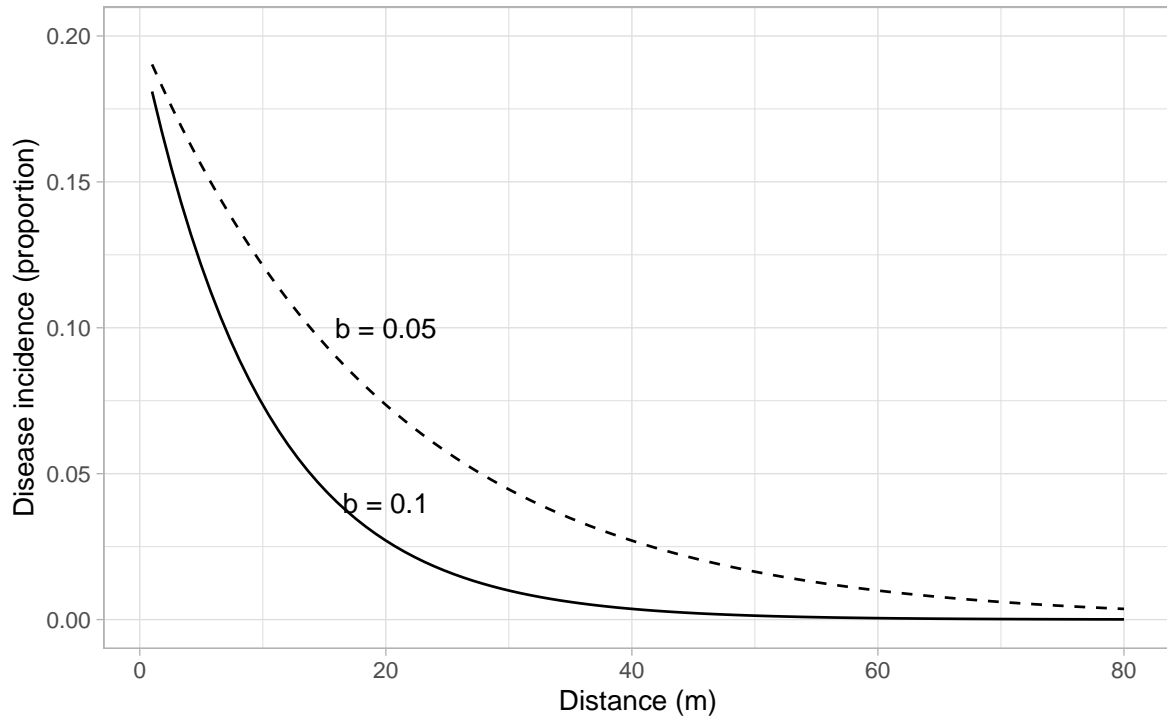
Set the parameters for the exponential model with two rates and same inoculum at the source:

```
a1 <- 0.2 # y at distance zero for gradient 1
a2 <- 0.2 # y at distance zero for gradient 2
b1 <- 0.1 # decline rate for gradient 1
b2 <- 0.05 # decline rate for gradient 2
max1 <- 80 # maximum distance for gradient 1
max2 <- 80 # maximum distance for gradient 2
dat <- data.frame(x = seq(1:max1), y = seq(0:a1))
```

The following code allows to visualize the model predictions.

```
dat %>%
  ggplot(aes(x, y)) +
  stat_function(fun = function(x) a1 * exp(-b1 * x), linetype = 1) +
  stat_function(fun = function(x) a2 * exp(-b2 * x), linetype = 2) +
  ylim(0, a1) +
  annotate("text", x = 20, y = 0.04, label = "b = 0.1") +
  annotate("text", x = 20, y = 0.10, label = "b = 0.05") +
  labs(
    title = "Exponential model",
    subtitle = "",
    x = "Distance (m)",
    y = "Disease incidence (proportion)"
  )
```

## Exponential model



### 2.1.2 Power law model

Also known as the modified Gregory's model (Gregory was a pioneer in the use this model to describe plant disease gradients). In the power law model,  $Y$  is proportional to the power of the distance, and is given by:

$$Y = a_P \cdot x - b_P$$

where  $a_P$  and  $b_P$  are the two parameters of the power law model. They differ from the exponential because as closer to  $x$  is to zero,  $Y$  is indefinitely large (not meaningful biologically). However, the model can still be useful because it produces realistic values at any distance  $x$  away from the source. The values of the  $a_P$  parameter should be interpreted in accord to the scale of  $x$ , whether in centimeters or meters. If the distance between the source and the first measure away from the source is 0.5m, it is so more appropriate to record the distance in cm than in m or km.

Once  $y$  at the distance zero from the source is undefined when using the power law model, this is usually modified by the addition of a positive constant  $C$  in  $x$ :

$$Y = a_P \cdot (x + C) - b_P$$

For this reason, the model is named as the modified power law. Here, the constant  $C$  is of the same unit of  $x$ . At the distance zero, the positive constant is a term that express the size of the inoculum source. In other words, the  $a$  parameter is a theoretical value of  $Y$  at the distance  $1 - C$  from the center of the inoculum source.

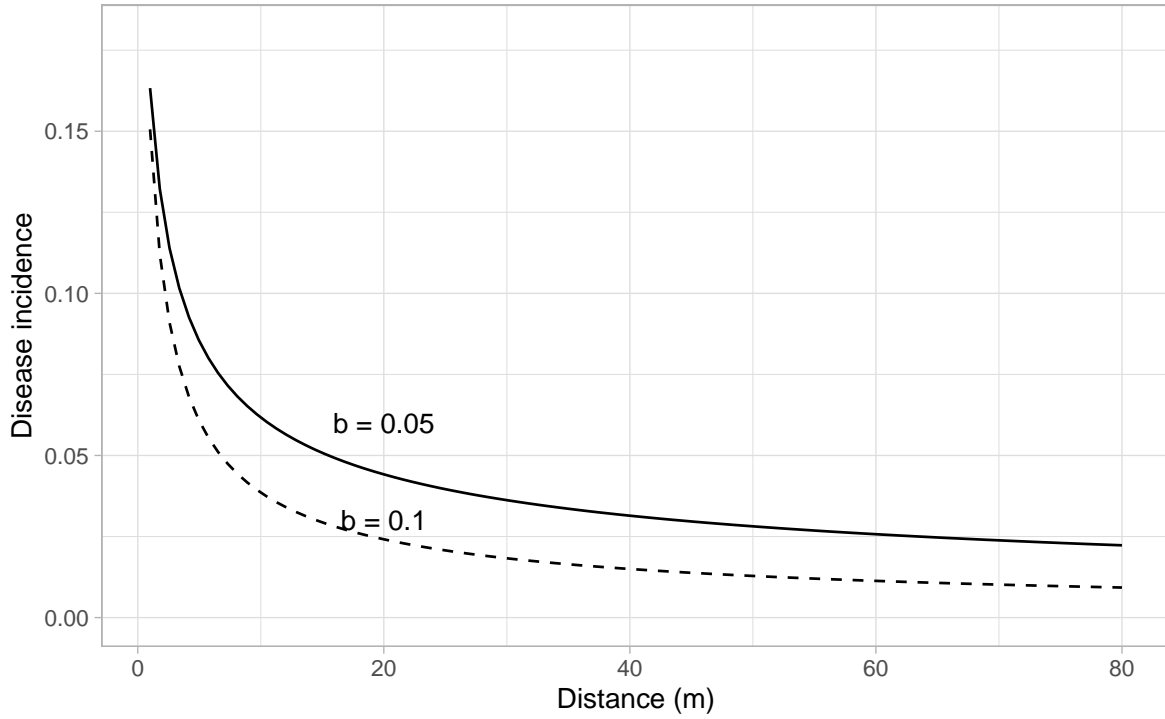
Let's plot two gradients with two rate parameters for the modified power law model:

```
C <- 0.5
a1 <- 0.2 # y at zero distance for gradient 1
a2 <- 0.2 # y at zero distance for gradient 2
b1 <- 0.5 # decline rate for gradient 1
b2 <- 0.7 # decline rate for gradient 2
max1 <- 80 # maximum distance for gradient 1
max2 <- 80 # maximum distance for gradient 2
dat2 <- data.frame(x = seq(1:max1), y = seq(0:a1))

dat2 %>%
  ggplot(aes(x, y)) +
  stat_function(fun = function(x) a1 * ((x + C)^-b1), linetype = 1) +
  stat_function(fun = function(x) a2 * ((x + C)^-b2), linetype = 2) +
  ylim(0, a1 - 0.02) +
  annotate("text", x = 20, y = 0.03, label = "b = 0.1") +
  annotate("text", x = 20, y = 0.06, label = "b = 0.05") +
  labs(
    title = "Modified Power Law",
    subtitle = "",
    x = "Distance (m)",
    y = "Disease incidence"
  )
```



## Modified Power Law



The differential equation of the power law model is given by:

$$\frac{dy}{dx} = \frac{-b_P \cdot Y}{x-C}$$

Similar to the exponential model,  $\frac{dy}{dx}$  is proportional to  $Y$ , meaning that the gradient is steeper (more negative) at the highest disease intensity value, usually closer to the source.

## 2.2 Linearization of the models

### 2.2.1 Transformations of $y$

The gradient models, again similar to the temporal disease models, are **non linear in their parameters**. The model is intrinsically linear if transformations are applied (according to the model) in both sides of the equations. The linear model in its generic state is given by

$$y^* = a^* + bx^*$$

where the asterisk in  $a$  indicated that one of the transformations was applied in  $y$  that produced the linear model. Note that  $a^*$  is the transformed version of the initial disease intensity, which

needs to be returned to the original scale according to the respective back-transformation. Follows the linearized form of the two most common gradient models.

$$\ln(y) = \ln(a_E) - b_E \cdot x$$

$$\ln(y) = \ln(a_P) - b_E \cdot \ln(x + C)$$

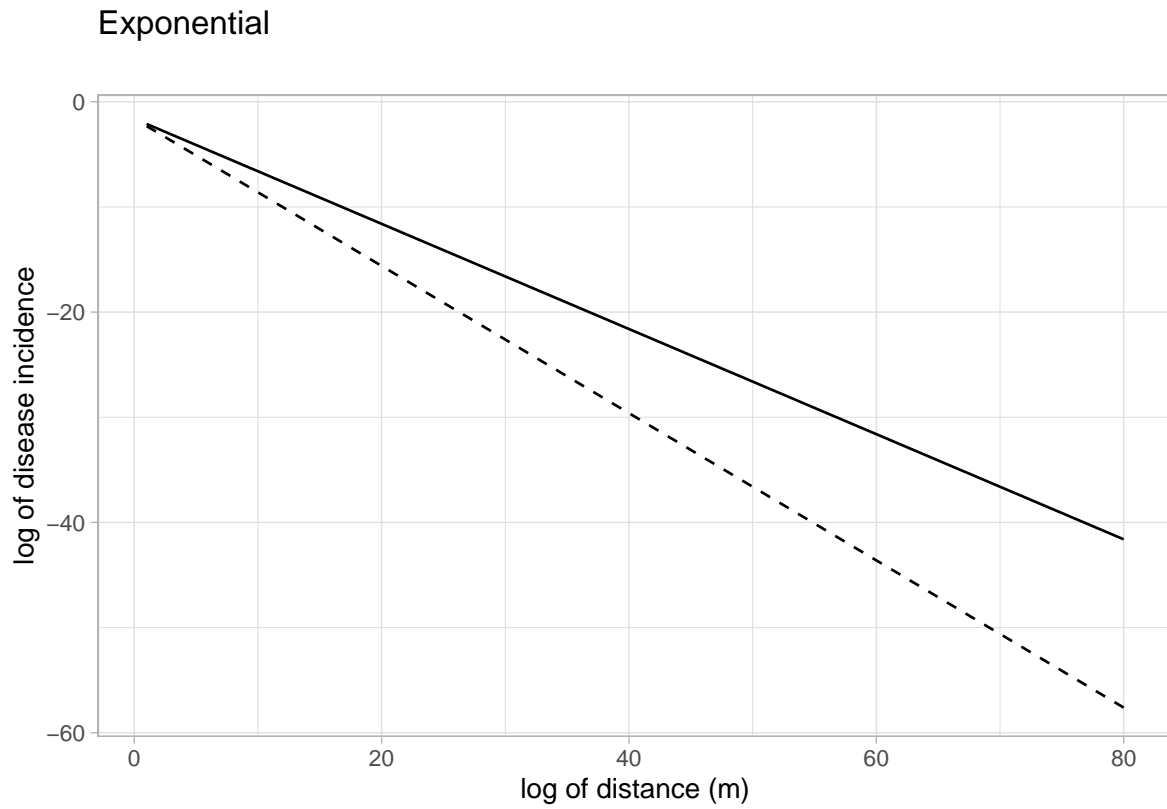
## 2.2.2 Plot for the linearized form of models

Let's visualize the linearization of the exponential model with two different slopes (gradient 1 and 2). Note that the transformation used was  $\ln(y)$ .

Follows the linearization of the modified power law model.

```
C <- 0.5
a1 <- 0.2 # y at zero distance for gradient 1
a2 <- 0.2 # y at zero distance for gradient 2
b1 <- 0.5 # decline rate for gradient 1
b2 <- 0.7 # decline rate for gradient 2
max1 <- 80 # maximum distance for gradient 1
max2 <- 80 # maximum distance for gradient 2
dat2 <- data.frame(x = seq(1:max1), y = seq(0:a1))

dat2 %>%
  ggplot(aes(x, y)) +
  stat_function(fun = function(x) log(a1) - (b1 * x), linetype = 1) +
  stat_function(fun = function(x) log(a2) - (b2 * x), linetype = 2) +
  labs(
    title = "Exponential",
    subtitle = "",
    x = "log of distance (m)",
    y = "log of disease incidence"
  )
```



Follows the linearization of the modified power law model. Note that the transformation used was  $\ln(y)$  and  $\ln(x + C)$  .

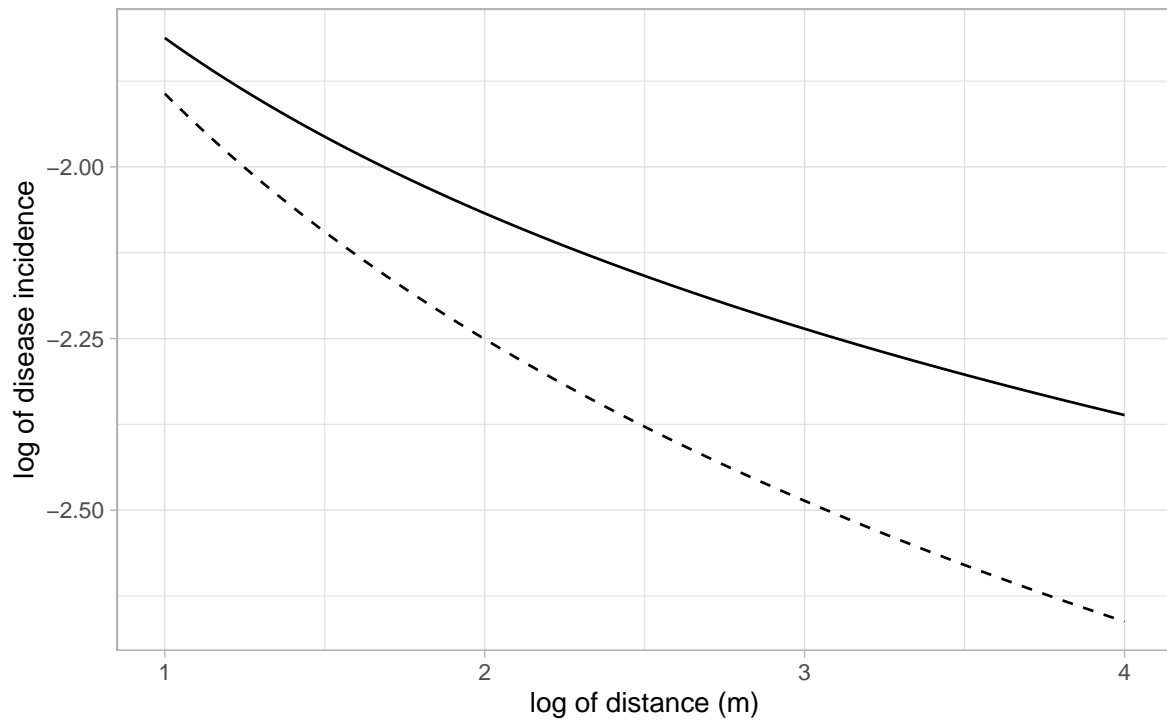
```

C <- 0.5
a1 <- 0.2 # y at zero distance for gradient 1
a2 <- 0.2 # y at zero distance for gradient 2
b1 <- 0.5 # decline rate for gradient 1
b2 <- 0.7 # decline rate for gradient 2
max1 <- log(80) # maximum distance for gradient 1
max2 <- log(80) # maximum distance for gradient 2
dat2 <- data.frame(x = seq(1:max1), y = seq(0:a1))

dat2 %>%
  ggplot(aes(x, y)) +
  stat_function(fun = function(x) log(a1) - (b1 * log(x + C)), linetype = 1) +
  stat_function(fun = function(x) log(a2) - (b2 * log(x + C)), linetype = 2) +
  labs(
    title = "Modified Power Law",
    subtitle = "",
    x = "log of distance (m)",
    y = "log of disease incidence"
  )

```

## Modified Power Law



## 2.3 Model fitting

### 2.3.1 Dataset

The hypothetical data below shows a gradient for the number of lesions counted at varying distances in meters from the source. Let's create two vectors, one for the distances  $x$  and the other for the lesion count  $Y$ , and then a dataframe by combining the two vectors.

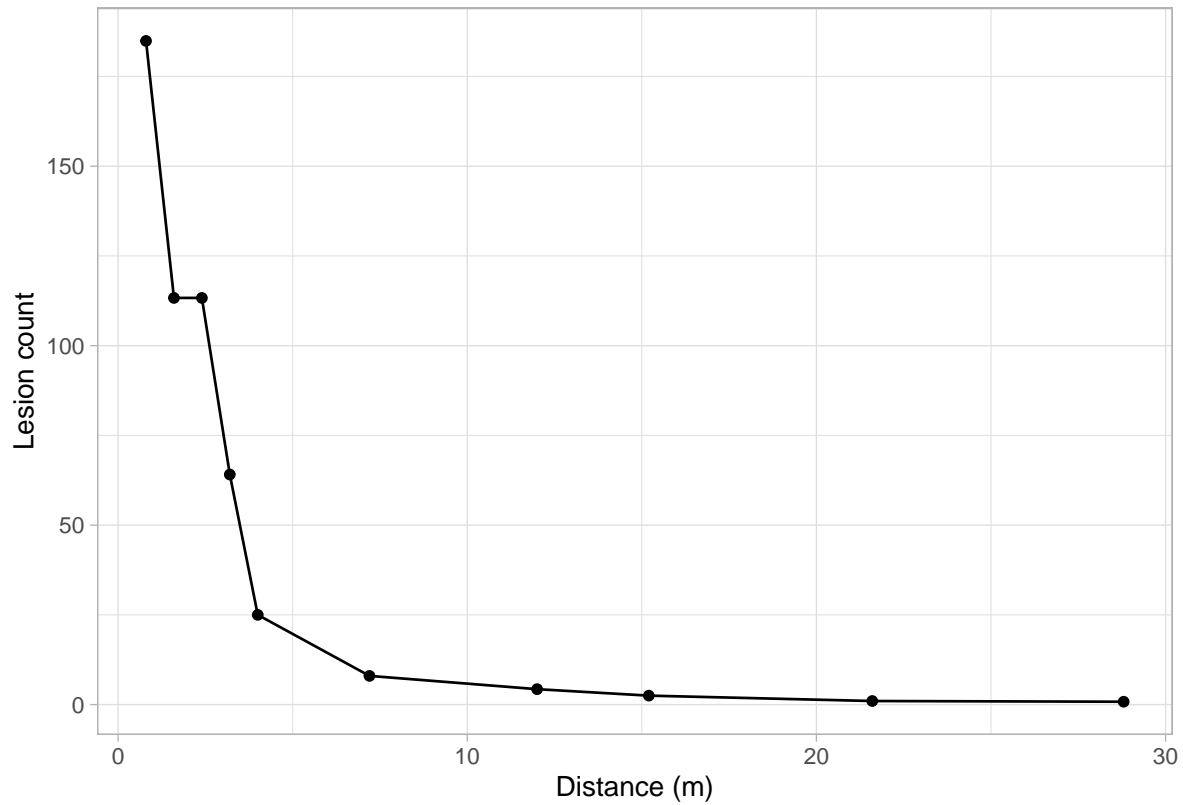
```
# create the two vectors
x <- c(0.8, 1.6, 2.4, 3.2, 4, 7.2, 12, 15.2, 21.6, 28.8)
Y <- c(184.9, 113.3, 113.3, 64.1, 25, 8, 4.3, 2.5, 1, 0.8)
grad1 <- data.frame(x, Y) # create the dataframe
grad1 # show the gradient
```

	x	Y
1	0.8	184.9

2	1.6	113.3
3	2.4	113.3
4	3.2	64.1
5	4.0	25.0
6	7.2	8.0
7	12.0	4.3
8	15.2	2.5
9	21.6	1.0
10	28.8	0.8

### 2.3.2 Visualize the gradient

```
grad1 %>%  
  ggplot(aes(x, Y))+  
  geom_point()+  
  geom_line()+  
  labs(y = "Lesion count",  
       x = "Distance (m)")
```



### 2.3.3 Linear regression

A linear regression model is fitted to the transformed variables according to the model. The higher the coefficient of determination, the better is the fit of the model to the data.

Exponential model

```
reg_exp <- lm(log(Y) ~ x, data = grad1)
summary(reg_exp)
```

Call:

```
lm(formula = log(Y) ~ x, data = grad1)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.04868	-0.58973	-0.00144	0.59572	0.99554

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	4.57705	0.35222	12.995	1.17e-06 ***
x	-0.20124	0.02656	-7.576	6.45e-05 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.7612 on 8 degrees of freedom

Multiple R-squared: 0.8777, Adjusted R-squared: 0.8624

F-statistic: 57.39 on 1 and 8 DF, p-value: 6.45e-05

Power law model with  $C = 0$ .

```
reg_p <- lm(log(Y) ~ log(x), data = grad1)
summary(reg_p)
```

Call:

```
lm(formula = log(Y) ~ log(x), data = grad1)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.72281	-0.11989	-0.03146	0.08755	0.65267

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.5638	0.2456	22.66	1.53e-08 ***
log(x)	-1.6978	0.1191	-14.26	5.71e-07 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4235 on 8 degrees of freedom

Multiple R-squared: 0.9621, Adjusted R-squared: 0.9574

F-statistic: 203.3 on 1 and 8 DF, p-value: 5.71e-07

Power law model with  $C = 0.4$ .

```
reg_pm <- lm(log(Y) ~ log(x + 0.4), data = grad1)
summary(reg_pm)
```



Call:

```
lm(formula = log(Y) ~ log(x + 0.4), data = grad1)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.53733	-0.17258	-0.03646	0.08450	0.56928

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.1007	0.2283	26.73	4.13e-09 ***
log(x + 0.4)	-1.8841	0.1084	-17.38	1.22e-07 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3495 on 8 degrees of freedom

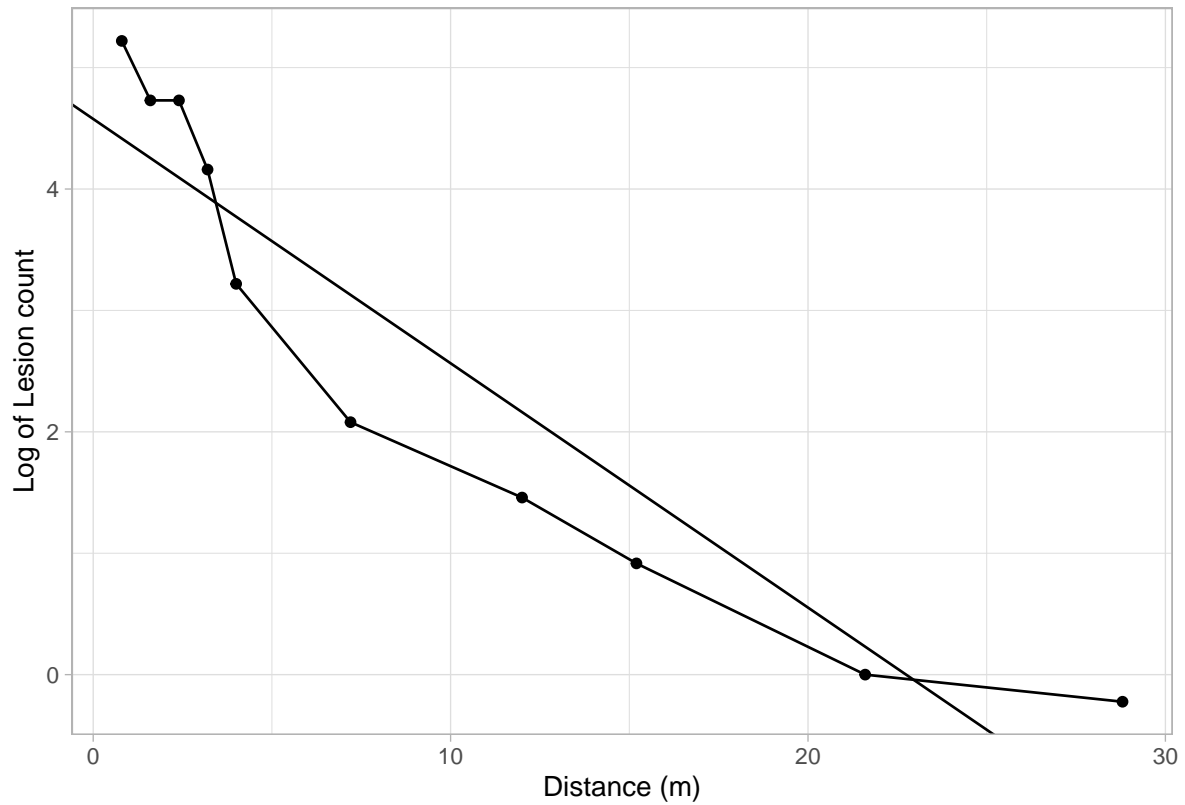
Multiple R-squared: 0.9742, Adjusted R-squared: 0.971

F-statistic: 302.2 on 1 and 8 DF, p-value: 1.223e-07

Graphs for the fitted models

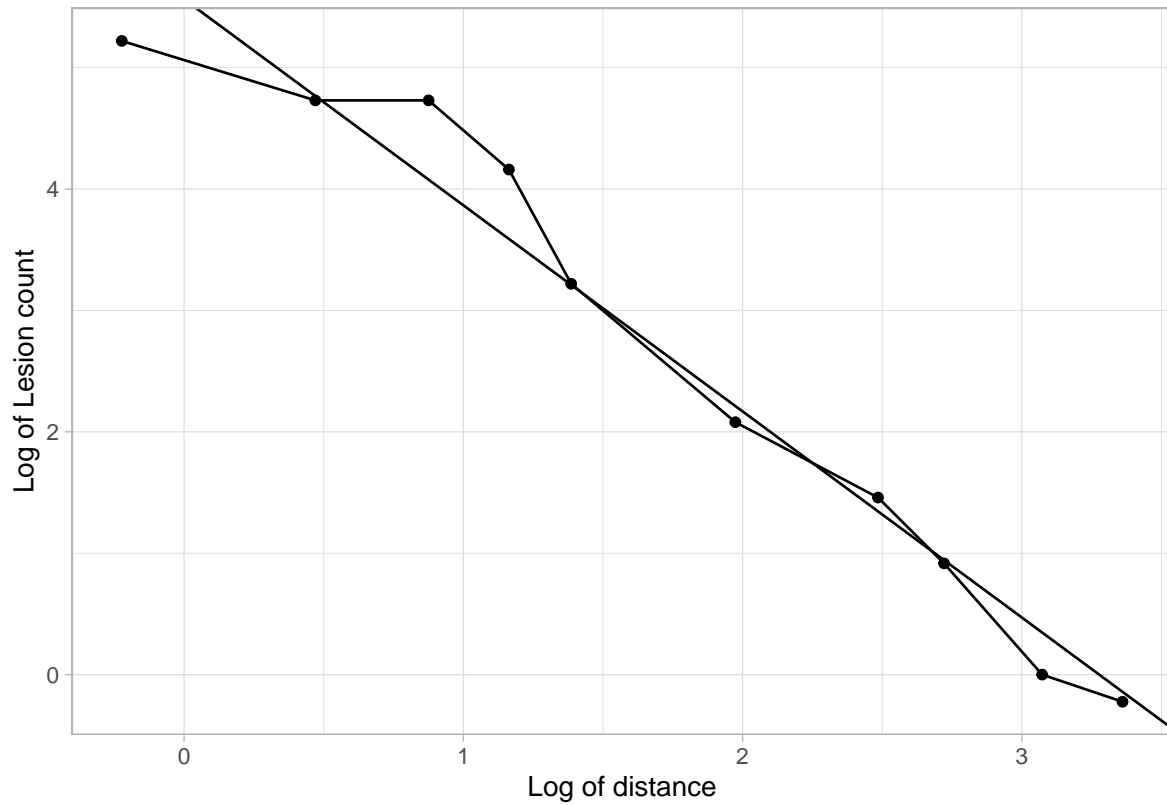
Exponential

```
grad1 %>%  
  ggplot(aes(x, log(Y)))+  
  geom_point()+  
  geom_line()+  
  geom_abline(slope = coef(reg_exp)[[2]], intercept = coef(reg_exp)[[1]])+  
  labs(y = "Log of Lesion count",  
       x = "Distance (m)")
```



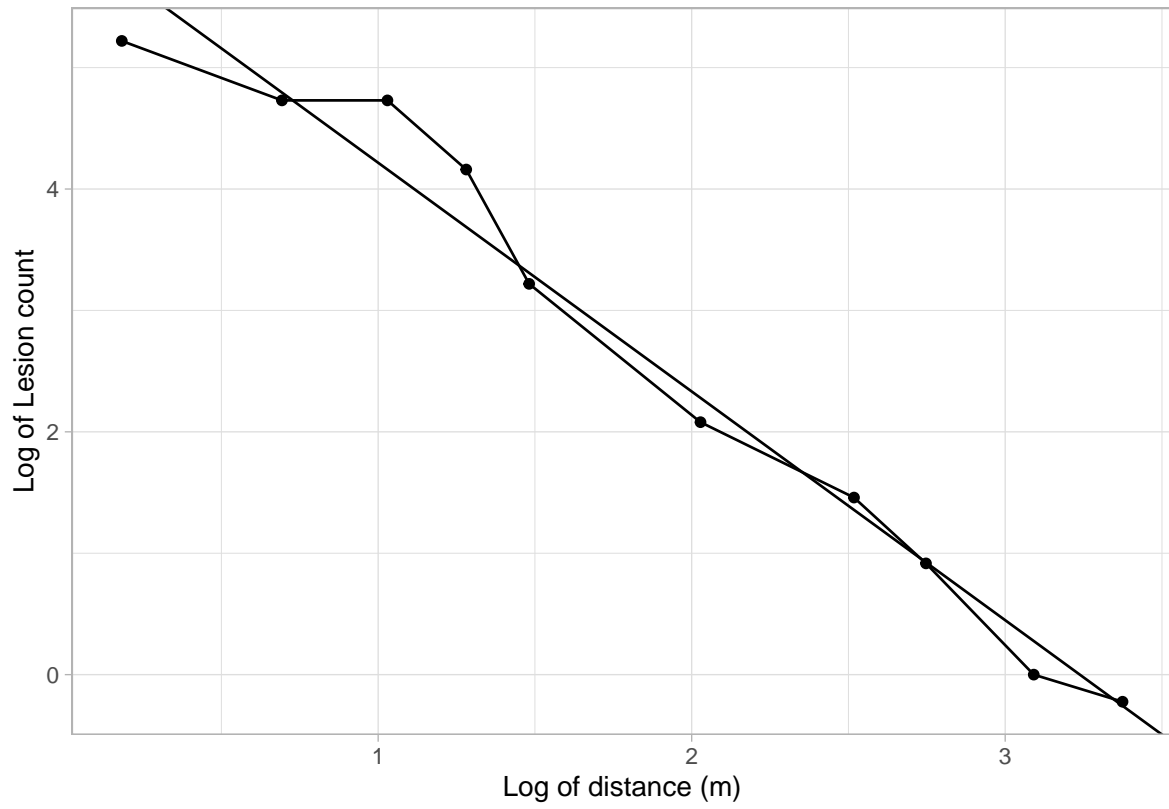
Power law model

```
grad1 %>%
  ggplot(aes(log(x), log(Y)))+
  geom_point()+
  geom_line()+
  geom_abline(slope = coef(reg_p)[[2]], intercept = coef(reg_p)[[1]])+
  labs(y = "Log of Lesion count",
       x = "Log of distance")
```



Modified power law model

```
grad1 %>%
  ggplot(aes(log(x+0.4), log(Y)))+
  geom_point()+
  geom_line()+
  geom_abline(slope = coef(reg_pm)[[2]], intercept = coef(reg_pm)[[1]])+
  labs(y = "Log of Lesion count",
       x = "Log of distance (m)")
```



Conclusion: The modified power law model provided the best fit.

## 3 Spatial patterns

### 3.1 Introduction

A spatial disease pattern can be defined as the arrangement of diseased entities relative to each other and to the architecture of the host crop (Madden et al. 2017b) . Such arrangement is the realization of the underlying dispersal of the pathogen, from one or several sources within and/or outside the area of interest, under the influence of physical, biological and environmental factors.

The study of spatial patterns is conducted at a specific time or multiple times during the epidemic. When assessed multiple times, both spatial and temporal processes can be characterized. Because epidemics change over time, it is expected that spatial patterns are not constant but change over time as well. Usually, plant pathologists are interested in determining spatial patterns at one or various spatial scales, depending on the objective of the study. The scale of interest may be a leaf or root, plant, field, municipality, state, country or even intercontinental area. The diseased units observed may vary from lesions on a single leaf to diseased fields in a large production region.

The patterns can be classified into two main types that occur naturally: **random** or **aggregated**. The random pattern originates because the chances for the units (leaf, plant, crop) to be infected are equal and low, and are largely independent from each other. In aggregated spatial patterns, such chances are unequal and there is dependency among the units, for example, a healthy unit close to a diseased unit is at higher risk than more distant units.

A range of techniques, most based on statistical tests, can be used to detect deviations from randomness in space and the choice of the methods depends on the scale of observation. Usually, more than one test is applied for the same or different scales of interest depending on how the data are collected. Three general categories of statistical tests can be determined based on the spatial scale and type of data collected: position of diseased units within a row or series of rows (plant to plant); quadrat or plot count data; or distance among the diseased units.

## 3.2 Position and status within a row (plant to plant)

Here, the status of each unit (usually a plant) is noted as a nominal variable. The plant is either diseased (D or 1) or non-diseased or healthy (H or 0). These data are usually collected within a crop row, giving rise to a series of binary data. Several statistical tests can be used to detect a deviation from randomness. The most commonly used tests are **runs**, **doublets** and **joint count** statistics.

### 3.2.1 Runs test

The runs test. Description.

Let's create a vector of binary (0 = non-diseased; 1 = diseased) data representing a crop row with 20 plants and assign it to `y`. For plotting purposes, we make a dataframe for more complete information.

```
y1 <- c(1,1,1,0,0,0,0,0,1,0,0,0,0,1,1,0,0,0,1,1)
x1 <- c(1:20) # position of each plant
z1 <- 1
row1 <- data.frame(x1, y1, z1) # create a dataframe
```

We can then visualize the series using ggplot.

```
library(tidyverse)
row1 %>%
  ggplot(aes(x1, z1, label = x1, color = factor(y1)))+
  geom_point(shape = 15, size = 6)+
  theme_void()+
  scale_x_continuous(breaks = max(z1))+
  scale_color_manual(values = c("green", "red"))+
  geom_text(vjust = 0, nudge_y = 0.5)+
  coord_fixed()+
  ylim(-0.5, 2.5)+
  theme(legend.position = "right")+
  labs(color = "Status", title = "Sequence of diseased (1) or non-diseased (0) units (plant)",
       subtitle = "The numbers represent the position of the unit")
```

Sequence of diseased (1) or non-diseased (0) units (plants)

The numbers represent the position of the unit



We can write a code in R and create a function named `oruns.test` for the ordinary runs test.

```

oruns.test <- function(x) {
# identify the sequence
S <- x
# Compute the number or runs
U = max(cumsum(c(1, diff(S)!=0)))
# Compute the number of diseased plants
m = sum(S)
# Count the total number of plants
N = length(S)
# Calculate the number of expected runs
EU = 1 + (2 * m*(N - m)/N)
# Calculate the standard deviation in the sample
sU = sqrt(2 * m * (N - m) * (2 * m * (N-m)-N) / (N^2 * (N-1)))
# Calculate the z-value
Z = (U - EU)/sU
# Obtain the p-value for the Z
pvalue <- (2*pnorm(abs(Z), lower.tail=FALSE))
# test if Z is lower than 1.64
result <- ifelse(Z < 1.64,
c("clustering"),
c("randomness"))
# Print the results
print(paste("There are",U,"runs. The number of expected runs is", round(EU,1), "P-value:",
})

```

We can now run the test for the example series above.

```

oruns.test(row1$y1)

```

```

[1] "There are 7 runs. The number of expected runs is 10.6 P-value: 0.084166 . Alternative hypothesis: clustering"

```

There are built-in functions in R packages that allow for running the ordinary runs test. Let's load the packages and run the test. Note that the results of the `runs.test` is the same as the one produced by our custom function.

```

library(randtests)
runs.test(row1$y1, threshold = 0.5)

```

Runs Test



```
data: row1$y1
statistic = -1.727, runs = 7, n1 = 8, n2 = 12, n = 20, p-value =
0.08417
alternative hypothesis: nonrandomness
```

```
library(DescTools)
r <- RunsTest(row1$y1)
```

### 3.2.2 Doublets

The doublets test. Description.

Let's manually produce a code to execute the doublets test. To facilitate, we can create a function and name it `doublets.test`. The only argument needed is the vector of binary data.

```

doublets.test <- function(x) {
# Identify the sequence
S <- x
# Compute the number of doublets Db
matrix <- cbind(S[-length(S)], S[-1])
pairs <- table(data.frame(matrix))
Db <- pairs[2,2]
# Count the number of diseased plants
N <- length(S)
# Count the number of total plants
m = sum(S)
# Expected number of doublets
EDb = m * ((m - 1) / N)
# Standard deviation
SDB = sqrt ( EDb * (1 - (2 / N)))
# Calculate the Z-value
ZDb = (Db - EDb) / SDB
# two-sided P-value calculation
pvalue <- (2 * pnorm(abs(ZDb), lower.tail = FALSE))
# Result of the test
result <- ifelse(abs(ZDb) >= 1.64,
c("aggregation or clustering"),
c("randomness"))
# Print the results
print(paste("There are", Db, "doublets. The number of expected doublets is", EDb, ".", "P-value: ", pvalue))
}

# Run the function calling the vector
doublets.test(row1$y1)

```

```
[1] "There are 4 doublets. The number of expected doublets is 2.8 . P-value: 0.4497 . Altern
```

### 3.2.3 Join count

Join count statistics. Description.

Let's use the `join.count` function of the `spdep` package to perform a joint count test. First we need to create the series of binary data from left to right. In the example, there are 5 rows and 5 columns. This will be informed later to run the test.

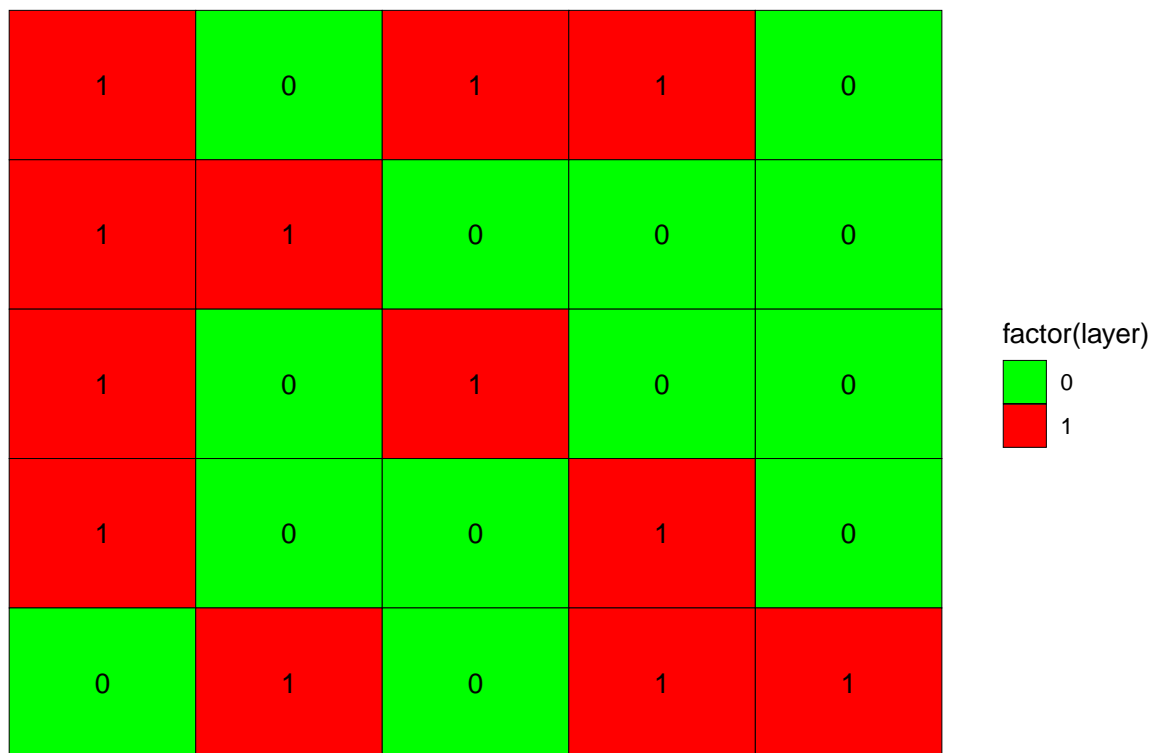
```
# Enter the data
S2 <- c(1,1,1,1,0,
        0,1,0,0,1,
        1,0,1,0,0,
        1,0,0,1,1,
        0,0,0,0,1)
```

Visualize the grid

```
# Convert to raster
mapS2 <- raster::raster(matrix(S2, 5 ,5))

# Convert to data frame
mapS3 <- raster::as.data.frame(mapS2,xy=TRUE)

# Map using ggplot
mapS3 %>%
  ggplot(aes(x, y, label = layer, fill = factor(layer)))+
  geom_tile(color = "black")+
  theme_void()+
  geom_text()+
  scale_fill_manual(values = c("green", "red"))
```



Load the library

```
library(spdep)
```

First, we need to generate a list of neighbors (nb) for a grid of cells. This is performed with the `cell2nb` function by informing the number of rows and columns. The argument “rook” means shared edge, but it could be the “queen”, for shared edge or vertex. We can use the default.

```
nb <- cell2nb(nrow = 5,
              ncol = 5,
              type="rook")
```

The `joincount.test` function runs the BB join count test for spatial autocorrelation. From the function description, the method uses a spatial weights matrix in weights list form for testing whether same-status joins occur more frequently than would be expected if the zones were labelled in a spatially random way. We need to inform the sequence as factor and the nb object we created previously.

```
joincount.test(factor(S2),
               nb2listw(nb))
```

Join count test under nonfree sampling

```
data: factor(S2)
weights: nb2listw(nb)
```

Std. deviate for 0 = -0.58266, p-value = 0.7199

alternative hypothesis: greater

sample estimates:

Same colour statistic	Expectation	Variance
2.9583333	3.2500000	0.2505797

Join count test under nonfree sampling

```
data: factor(S2)
weights: nb2listw(nb)
```

Std. deviate for 1 = -0.66841, p-value = 0.7481

alternative hypothesis: greater

sample estimates:

Same colour statistic	Expectation	Variance
2.4166667	2.7500000	0.2486957

The function returns a list with a class for each of the status (in this case 0 and 1) with several components. We should look at the **P-value**. The alternative hypothesis (greater) is that the same status joins occur more frequently than expected if they were labelled in a spatial random way. In this case, we do not reject the null hypothesis of randomness.

We can run the ordinary runs and doublets tests, which only considers the adjacent neighbor, for the same series and compare the results.

```
oruns.test(S2)
```

```
[1] "There are 13 runs. The number of expected runs is 13.5 P-value: 0.844252 . Alternative 1"
```

```
doublets.test(S2)
```

[1] "There are 5 doublets. The number of expected doublets is 5.28 . P-value: 0.8989 . Altern

## References

- Madden, L. V., Hughes, G., and van den Bosch, F., eds. 2017a. CHAPTER 4: Temporal analysis i: Quantifying and comparing epidemics. In *The American Phytopathological Society*, p. 63–116. Available at: <http://dx.doi.org/10.1094/9780890545058.004>.
- Madden, L. V., Hughes, G., and van den Bosch, F., eds. 2017b. CHAPTER 9: Spatial aspects of epidemics—III: Patterns of plant disease. In *The American Phytopathological Society*, p. 235–278. Available at: <http://dx.doi.org/10.1094/9780890545058.009>.