

# **Mini Project :Population Growth**

**Liyang Huang**  
**supervisor: Dr Samraat Pawar**  
**Faculty of Natural Sciences,,**  
**Department of Life Sciences (Silwood Park),**  
**s.pawar@imperial.ac.uk**

**March 2020,**  
**WORD COUNT: 2563**

# 1 Introduction

Natural populations are composed of individuals with diverse phenotypes that show differences in their demographic parameters and intra- and inter-species interaction([1]. Fluctuations in individual population abundance may play a key role in ecosystem dynamics and emerging functional characteristics. When the abundance is low and the resources are not limited, the population abundance will increase exponentially, which is the Malthus principle ([3]. The Malthusian principle also points out that when resources become limited, population growth will gradually slow down and eventually stop. At the same time, there may be a period of time before population growth really starts, which is called the lag phase. The data used in this simulation was collected through laboratory experiments around the world, ensuring the sample size of this experiment. In this report, I used different models to simulate the growth process of population abundance, and made comparison and analysis to obtain the most suitable model for the growth process of experimental samples.

## 2 Background

This project using 4 model, the logistic equation model, the modified Gompertz model, The Baranyi model and The Buchanan model.

### 2.1 The logistic equation model

This model is the simplest mathematical models that we can use the phenomenological quadratic and cubic polynomial models. In general, if the quantitative characteristics of objective things are: at time  $t$  is small, things grow exponentially, and when  $t$  increases, the growth rate gradually decreases and gets closer and closer to a certain value (that is, bearing capacity  $N_{max}$ ), such problems can be solved by the Logistic equation ([4].

$$N_t = \frac{N_0 N_{max} e^{rt}}{N_{max} + N_0(e^{rt} - 1)}$$

27

28  $N_t$  is population size at time  $t$

29  $N_0$  is initial cell culture (Population) density

30  $N_{max}$  is maximum population density, called carrying capacity

31  $t$  is time that is parameters in the sample data

32  $r$  is maximum growth rate

## 33 2.2 The modified Gompertz model

34 This model is often used in the literature to simulate bacterial growth. It  
 35 is a function of the sigmoid colon and describes the slowest growth at the  
 36 beginning and end of a given time period. This is the most widely accepted  
 37 detailed convention on population growth. The right-hand or future value  
 38 asymptote(A) of the function is closer to the curve than the left Side or  
 39 lower asymptotes ([6]).

$$N_t = Ae^{-e^{\frac{r_{max}(t_{lag}-t)}{A}+1}}$$

40  
 41  $N_t$  is population size at time  $t$   
 42  $N_0$  is initial cell culture (Population) density  
 43  $N_{max}$  is maximum population density, called carrying capacity  
 44  $t$  is time that is parameters in the sample data  
 45  $r_{max}$  is maximum growth rate  
 46  $t_{lag}$  is the x-axis intercept to this tangent, means duration of the delay  
 47 before the population starts growing exponentially  
 48  $A = \ln(N_{max}/N_0)$ , is the asymptote.

## 49 2.3 The Baranyi model

50 Four kinds of common survival curves are fitted with Baranyi model: linear  
 51 curve, hysteresis curve, trailing curve and sigmoid curve. This model adds a  
 52 new dimensionless parameter  $h_0$  which represents the initial physiological  
 53 state of the cells. For the prediction performance, the Baranyi model was  
 54 better and more robust than the modified Gompertz equation ([5]).

$$N_t = N_0 + r_{max}A_t - \ln\left(1 + \frac{e^{r_{max}A_t} - 1}{e^{N_{max}N_0}}\right),$$

55  
 56 where:

$$A_t = t + \frac{1}{r_{max}} \cdot \ln\left(\frac{e^{-r_{max}t} + h_0}{1 + h_0}\right).$$

$$t_{lag} = \frac{\ln\left(1 + \frac{1}{h_0}\right)}{r_{max}},$$

57

58  $N_t$  is population size at time  $t$

59  $N_0$  is initial cell culture (Population) density

60  $N_{max}$  is maximum population density, called carrying capacity

61  $t$  is time that is parameters in the sample data

62  $r_{max}$  is maximum growth rate

63  $t_{lag}$  is the x-axis intercept to this tangent, means duration of the delay

64 before the population starts growing exponentially

## 65 2.4 The Buchanan model

66 This model can be called as three-phase logistic model. Three-phase is an  
67 Initial Phase, Intermediate Phase, and Final Phase. The initial stage means  
68 that  $t_{lag}$ , relatively stable, or flat over time. The Intermediate Phase refers  
69 to  $t_{lag} \leq t \leq t_{max}$ . After the initial stage, the growth rate may change. If  
70 the initial population is much smaller than the carrying capacity, the pop-  
71 ulation will increase rapidly. If the initial population abundance is much  
72 larger than the carrying capacity, the population will decrease rapidly. If  
73 the initial population abundance approaches the capacity, the population  
74 abundance will tend to be stable. Final phase means  $t_{lag} \leq t$ , when the pop-  
75 ulation abundance reach the carrying capacity. At this point, the population  
76 abundance will be stable, unless the carrying capacity changes([2].

$$N(t) = \begin{cases} N_0 & \text{if } t \leq t_{lag} \\ N_{max} + r_{max} \cdot (t - t_{lag}) & \text{if } t_{lag} < t < t_{max} \\ N_{max} & \text{if } t \geq t_{max} \end{cases}$$

77

78  $N_t$  is population size at time  $t$

79 N0 is initial cell culture (Population) density  
80 Nmax is maximum population density, called carrying capacity  
81 t is time that is parameters in the sample data  
82 rmax is maximum growth rate  
83 tlag is the x-axis intercept to this tangent, means duration of the delay  
84 before the population starts growing exponentially

## 85 3 methods and data

### 86 3.1 Computing tools

87 In this project, I use three scripting language to write this project. Python,  
88 R and Script. Python is used to filter data from sample data and find  
89 some parameter that I will use in the next program. R is used to simulate  
90 the model and perform the nls operation (fit the model), while generating  
91 pictures and final data. Script is used to run all programs, reducing the  
92 time required to manually run the programs and increasing the efficiency of  
93 the project. In python, I call a lot of packages to reduce some unnecessary  
94 runtime and memory usage. I write two python program. The first program  
95 is used to filter data. For this purpose, I called pandas, seaborn, math as a  
96 package. Pandas used to be create new DataFrame to store data.

```
data = pd.read_csv("../data/LogisticGrowthData.csv")  
pd.read_csv("../data/LogisticGrowthMetaData.csv")
```

97  
98 I prefer to use DataFrames to organize data over lists. Using DataFrames  
99 to organize time faster, reduce runtime and memory data.

100

101 Seaborn used to plot.

```
#sns.lmplot("Time", "PopBio", data = data_subset, fit_reg = False)
```

102

103 In this project, using it will plot the abundance of the population over time.  
104 But this picture is only used for the parameter reference is not needed at  
105 the end, so comment out in the code

106

107 Math used to calculate  $\ln(N_{MAX} / N_0)$ . In python, `math.log()` means  
108 return the natural logarithm of x.

```
LogPopBio = math.log(PopBio.iat[i,0])
```

109

110 The second program is used to find some parameter that I will use in the  
111 next program. For this purpose, I just called pandas. Using DataFrame to  
112 store data that I will use in the next program.

```
data = pd.read_csv("../data/data.csv")
ID = pd.read_csv("../data/ID.csv")
rmaxList = []
```

113

114 In r file, I called repr, ggplot2, nls.multstart as a package. Repr used to  
115 change default plot size.

```
library(nls.multstart)
options(repr.plot.width=4, repr.plot.height=4) # Change default p
require(["minpack.lm"])
```

116

117 ggplot2 used to plot. In this project, using it can plot the abundance of the  
118 population over time and model and sample data matching graph.

```
g1<- ggplot(DF, aes(x = Time, y = LogPopBio)) +
  geom_point(size = 3) +
  geom_line(data = model_frame, aes(x = Time, y = LogPopB
  theme_bw(base_size = 16) + # make the background white
  theme(aspect.ratio=1)+ # make the plot square
  labs(x = "Time", y = "log(PopBio)")
  #print(g1)
```

119

120 nls.multstart used to NLLS fitting method. Using it to match the sample  
121 data to the four models respectively.

```

timepoints <- DF$Time
fit_logistic <- try(nlsLM(LogPopBio ~ logistic_r
| | | | | | | | | | list(rmax=Rmax, N0 = N0, N
if('try-error' %in% class(fit_logistic)){
  print("error")
  logistic_points <- rep(0,nrow(DF))
  df1 <- data.frame(timepoints, logistic_points)
  df1$model <- "Logistic"
  names(df1) <- c("Time", "LogPopBio", "model")
  AIC_logistic <- 0
  BIC_logistic <- 0

```

122  
 123 In this part, also using “try-error” method. This part will described in  
 124 ‘method’ section.

## 125 3.2 Data

126 In this project, all data comes from research institutes around the world,  
 127 ensuring that the amount of sample data is large enough. The sample data  
 128 set is called LogisticGrowthData.CSV. It contains measurements of change  
 129 in biomass or number of cells of microbes over time. Meanwhile, we will use  
 130 the other data set is called LogisticGrowthMetaData.CSV. This contains  
 131 a detailed description of each parameter. Using panda.DataFrame to read  
 132 these data set to our program. When we filter the data, all the filtered data  
 133 is written into a new CSV, called data.CSV. The ID name corresponding to  
 134 the filtered data is written to another CSV file, called ID.CSV. In all the  
 135 following processes, we only use data.CSV and ID.CSV to process the data.

## 136 3.3 Method

137 The propose of this project is compare the selected sample data with the  
 138 population growth calculated by the model and to obtain the model that  
 139 most conforms to the sample data.

### 141 3.3.1 Process Data

142 For this purpose, we first set a unique ID for the sample data to group  
 143 the sample data. In this project, I chose to evaluate the uniqueness of  
 144 all parameters and make the parameter values part of the ID rather than  
 145 unique. Therefore, the final ID is data.Species + data.Temp.map (str) +  
 146 data.Medium + data.Citation + data.Rep.map ( str). The role of map is  
 147 to convert int to string, because in the id naming process, only parameters

with type of string can be accepted. After grouping the sample data, we got 305 sets of data. Next, according to the parameters to be used,  $r_{max}$ ,  $n_0$ ,  $n_{max}$ ,  $t_{lag}$ , and  $A$ , we extract each group with a total number of 5 or more into a single data group, and gather the corresponding ID names to build a data set. Next, we need to find the parameters use in the mathematical formula. The used parameters are  $n_0$ ,  $n_{max}$ ,  $r_{max}$ ,  $t_{lag}$  and  $a$ . Among them,  $r_{max}$  refers to the maximum slope. In order to find the maximum slope, we take the simplest approach,  $(K = y_1 - y_2 / x_1 - x_2)$ . First we set the first point as the starting point, and then use the first and second points to find the slope. After that we use a for loop to find the slope between each two points, and finally get the maximum slope. Find the value of  $k$  while recording the  $x$  and  $y$  values, and get the value of  $b$  by  $y = kx + b$ .  $n_0$  is the minimum value in the data,  $n_{max}$  is the maximum value in the data,  $t_{lag}$  is the value of  $x$  when  $y = n_0$ , and  $A = \ln(n_{max} / n_0)$ . After all the parameters of the first data group were found, the for loop was going to be used to find the parameters of the next data group, and the parameters obtained were combined with the parameters of the previous data group to form a new data frame. Finally, the generated data frames are written to a CSV file(LIST.CSV).

### 3.3.2 Model fitting

When we have processed all the data, we use R for model fitting. Read the LIST.CSV file and data.CSV file into the DataFrame, and use the subset to find the data corresponding to the same id. Next, build the model and get the parameters needed for the model. Use the NLLS fitting method to fit the model and actual data. The NLLS fitting method is used to fit the model and the actual data, the COEF value of the obtained model is extracted as a data frame, and calculated AIC and BIC. In the NLLS fitting method, some models can not be matched. Use the try-catch method to avoid errors. Finally, import the results from the previous step and print the model. Visually identify bad fits and determine whether to further optimize previous NLLS fit scripts. Meanwhile, AIC and BIC were imported to analyze the fitting results of the model and summarize the most suitable model.

## 4 Results

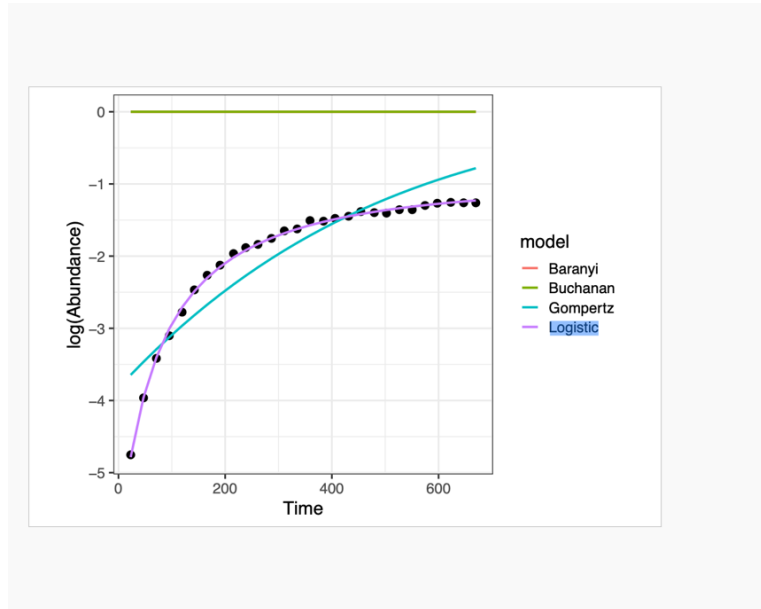
The propose of this project is to compare the selected sample data with the population growth calculated by the model, and to obtain the model that most conforms to the sample data. In this project, python, r, and script. For loop, if, try-catch, NLLS, DataFrame and other methods are used. The result contains 288 images and a CSV file, which contains AIC and BIC.



188 The following uses some fitting models to illustrate the results.  
 189

## 190 4.1 Logistic Model

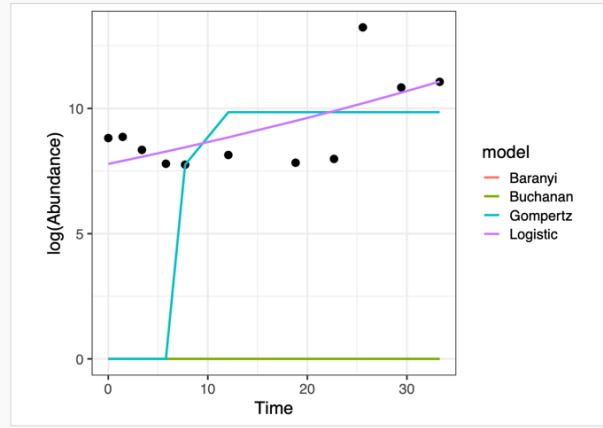
191 The data used in this case is id Chryseobacterium.balustinum-5-TSB-Bae,  
 192 YM, Zheng, L., Hyun, JE, Jung, KS, Heu, S. and Lee, SY, 2014. Growth  
 193 characteristics and biofilm formation of various spoilage bacteria isolated  
 194 from fresh produce. Journal of food science, 79 (10), pp.M2072-M2080.-  
 195 1,  $r_{max} = 0.000344284$ ,  $n_0 = -4.752591912$ ,  $n_{max} = -1.254736278$ , AIC-  
 196 gompertz= 29.93912607, AIC-logistic= -100.4740036, BIC-gompertz= 35.26794411,  
 197 BIC-logistic= -95.14518554



198  
 199 This figure shows that the Logistic model is the most suitable model. Baranyi  
 200 model and Buchanan are not suitable for this data set. The curve of logistic  
 201 model is generally J-shaped. As can be seen from this figure, it conforms to  
 202 this special feature. When  $t$  grows from small to large, the population abun-  
 203 dantness keeps increasing. However, when the population abundance reaches a  
 204 certain value, it reaches a stable state with the increase of  $t$ . From the char-  
 205 acteristics of the logical model, if the quantitative characteristics of objective  
 206 things are as follows: when time  $t$  is very small, things increase exponen-  
 207 tially; when  $t$  increases, the growth rate gradually decreases and gets closer  
 208 and closer to a certain value (that is, bearing capacity  $N_{max}$ ), the logical  
 209 model is most suitable for this data set.

## 210 4.2 Gompertz Model

211 The data used in this case is id *Lactobacillus sakei*-30-MRS broth-Silva,  
 212 A.P.R.D., Longhi, D.A., Dalcanton, F. and Arago, G.M.F.D., 2018. Mod-  
 213 elling the growth of lactic acid bacteria at different temperatures. Brazil-  
 214 ian Archives of Biology and Technology, 61.-1,  $r_{max} = 0.920853841$ ,  $n_0$   
 215  $= 8.817495339$ ,  $n_{max} = 13.23583879$ ,  $AIC-gompertz = 75.9798774$ ,  $AIC-$   
 216  $logistic = 45.20044445$ ,  $BIC-gompertz = 77.5714585$ ,  $BIC-logistic = 46.79202554$

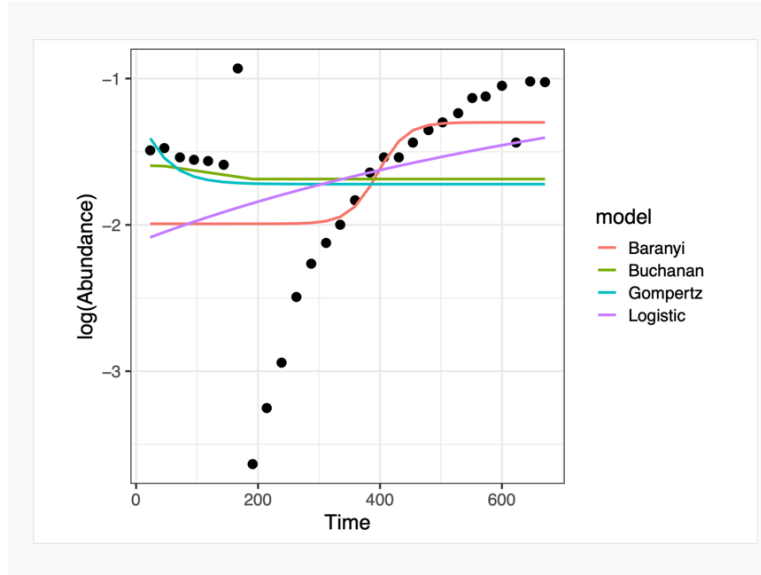


217 This figure shows that the Gompertz model is the most suitable model.  
 218 Baranyi model and Buchanan are not suitable for this data set. The logis-  
 219 tic model curve is generally s-shaped. As can be seen from this figure, it  
 220 conforms to this special feature. When  $t$  starts to increase, the growth rate  
 221 of population abundance slightly slower, and the growth rate of population  
 222 abundance is accelerated in the middle period, and is slower in the later  
 223 period. The Gompertz model is best suited to this data set in terms of the  
 224 characteristics of a logical model that grows slowest at the beginning and  
 225 end of a given time period.  
 226

## 227 4.3 Baranyi model

228 The data used in this case is id *Clavibacter michiganensis*-5-TSB-Bae, Y.M.,  
 229 Zheng, L., Hyun, J.E., Jung, K.S., Heu, S. and Lee, S.Y., 2014. Growth char-  
 230 acteristics and biofilm formation of various spoilage bacteria isolated from  
 231 fresh produce. Journal of food science, 79(10), pp.M2072-M2080.-1,  $r_{max}$   
 232  $= 0.015869066$ ,  $n_0 = -1.437433323$ ,  $n_{max} = -0.931062049$ ,  $AIC-gompertz =$

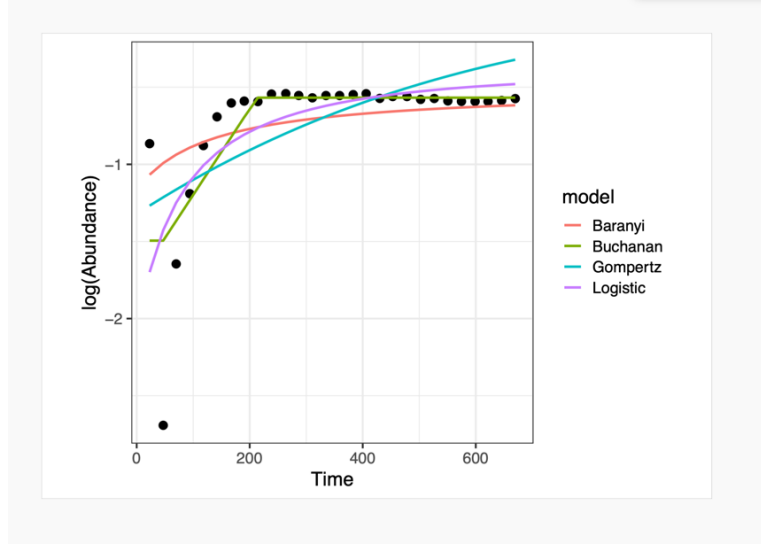
233 64.31968967,AIC-logistic= 60.76109259,AIC-baranyi= 55.28514297, AIC-buchanan=  
 234 66.32887698,BIC-gompertz= 69.64850771,BIC-logistic= 66.08991063, BIC-  
 235 baranyi= 61.94616552, BIC-buchanan= 66.32887698



236  
 237 This figure shows that the Baranyi model is the most suitable model. The  
 238 curves of baranyi model are generally linear, with hysteresis phase, trailing  
 239 phase and s-shaped curve. As can be seen from this figure, it conforms to  
 240 this special feature. When  $t$  starts to increase, the population abundance is  
 241 already at a high level, indicating that the cells have an initial physiological  
 242 state ( $h_0$ ). According to the characteristics of the logical model, a new  
 243 dimensionless parameter  $h_0$  is added to represent the initial physiological  
 244 state of the cell. Therefore, the Baranyi model is best suited for this data  
 245 set.

#### 246 4.4 Buchanan model

247 The data used in this case is id Enterobacter.sp.-5-TSB-Bae, Y.M., Zheng,  
 248 L., Hyun, J.E., Jung, K.S., Heu, S. and Lee, S.Y., 2014. Growth char-  
 249 acteristics and biofilm formation of various spoilage bacteria isolated from  
 250 fresh produce. Journal of food science, 79(10), pp.M2072-M2080.-1,  $r_{max}$   
 251 =  $7.50E-05$ ,  $n_0 = -0.865921996$ ,  $n_{max} = -0.541427051$ , AIC-gompertz=  
 252 30.20418465, AIC-logistic= 23.08963106, AIC-baranyi= 34.07912787, AIC-  
 253 buchanan= 16.8633038,BIC-gompertz= 35.53300269,BIC-logistic= 28.4184491,  
 254 BIC-baranyi= 40.74015042, BIC-buchanan= 16.8633038



255

256 This figure shows that the Buchanan model is the most suitable model. The  
 257 curve of Buchanan model is generally three-phase, and the three phases are  
 258 the Initial phase, the Intermediate phase, and the Final phase. As can be  
 259 seen from this figure, it conforms to this special feature. When  $t$  starts to  
 260 increase, the population abundance is relatively stable. Over time, the  
 261 initial population was less than the carrying capacity, and the population  
 262 abundance increased rapidly. As the population abundance approaches the  
 263 carrying capacity, the population abundance remains stable. From the char-  
 264 acteristics of the logical model, Initial phase means  $t \leq t_{lag}$ , and is relatively  
 265 stable or flat over time. Intermediate Phase means  $t_{lag} < t < t_{max}$ . After the  
 266 initial phase, growth rates may change. If the initial population is much  
 267 smaller than the carrying capacity, the population will increase rapidly. If  
 268 the initial population abundance is much greater than the carrying capac-  
 269 ity, the population population will decrease rapidly. If the initial population  
 270 abundance is close to the volume, the population abundance tends to be sta-  
 271 ble. Final phase means  $t_{lag} \leq t$  when the population abundance reach the  
 272 carrying capacity. At this point, the population abundance will be stable,  
 273 unless the carrying capacity changes ([2]. Therefore, the Buchanan model  
 274 is most suitable for this data set.

## 275 5 Discussion

276 This project verified the match between the three models and the data, but  
 277 not all the data can be 100 percent matched. The problem should be that  
 278 the data filter is not carefully classified.

## 279 References

- 280 [1] Daniel I Bolnick, Priyanga Amarasekare, Márcio S Araújo, Reinhard  
281 Bürger, Jonathan M Levine, Mark Novak, Volker HW Rudolf, Sebas-  
282 tian J Schreiber, Mark C Urban, and David A Vasseur. Why intraspe-  
283 cific trait variation matters in community ecology. *Trends in ecology &  
284 evolution*, 26(4):183–192, 2011.
- 285 [2] Sophie López, M Prieto, Jan Dijkstra, Mewa Singh Dhanoa, and Jim  
286 France. Statistical evaluation of mathematical models for microbial  
287 growth. *International journal of food microbiology*, 96(3):289–300, 2004.
- 288 [3] Thomas Robert Malthus, Donald Winch, and Patricia James.  
289 *Malthus: 'An Essay on the Principle of Population'*. Cambridge Uni-  
290 versity Press, 1992.
- 291 [4] Eric W Weisstein. Logistic equation. 2003.
- 292 [5] R Xiong, G Xie, AS Edmondson, RH Linton, and MA Sheard. Com-  
293 parison of the baranyi model with the modified gompertz equation for  
294 modelling thermal inactivation of listeria monocytogenes scott a. *Food  
295 Microbiology*, 16(3):269–279, 1999.
- 296 [6] MH Zwietering, Il Jongenburger, FM Rombouts, and KJAEM  
297 Van't Riet. Modeling of the bacterial growth curve. *Appl. Environ.  
298 Microbiol.*, 56(6):1875–1881, 1990.