

Comparisons of Models Applied in Microbial Population Growth

Congjia.chen (Congjia.Chen21@imperial.ac.uk)

Dec, 2021

¹Word Count: 3058

²CMEECoursework, Department of Life Sciences, Imperial College of Science, Technology and Medicine

Abstract

1 Microbial growths is highly correlated with human society. There-
2 fore, having knowledge of the microbial growth is essential so that
3 human can anticipate or control their growth under particular con-
4 ditions. Mathematical models have been proved to be functional in
5 microbial growth anticipation. However, lack of empirical model com-
6 parison with universal data will lead some bias to the model selec-
7 tion. In this report, based on the model fitting and model selection
8 on 285 empirical data sets, non-linear model performs better than lin-
9 ear model. Among the non-linear model, logistic model is sufficient for
10 simple growth situation, while Gompertz and Baranyi can handle more
11 complicated situation like lag phase. However, all of the non-linear
12 models such as logistic model and Baranyi model have the defect that
13 they can not be applied to fit the death phase properly. Segmented
14 model and external factors calibration might be a potential strategy
15 to optimize the result.

1 Introduction

16 Microbial growths is highly correlated with human society, for example, the
17 yeast growing in wort to make beer, the pathogenic bacteria to make hu-
18 man sick and the microbe which leads to food spoilage. Therefore, having
19 knowledge of the microbial growth is essential so that human can antic-
20 ipate or control their growth under particular conditions[1]. In contrast
21 to multi-cellular organisms, microbial growth is measured by population
22 growth, either by counting the number of cells or by increasing the overall
23 mass. However, current methods of population measurement are relatively
24 complicated[2].

25 In a closed system, the growth curve of the microbial population can

26 be divided into four phases: (1) Lag Phase, (2) Log(exponential) phase,
27 (3) Stationary phase, (4) Death phase[3]. In particular, after adapting the
28 new environment to the lag phase, the microbial population increases ex-
29 ponentially while its abundance is low and resources are not limited. This
30 growth slows down and ends up when resources become scarce.

31 Figure 1 shows that the application of mathematical model is dramati-
32 cally increasing on microbial researches.

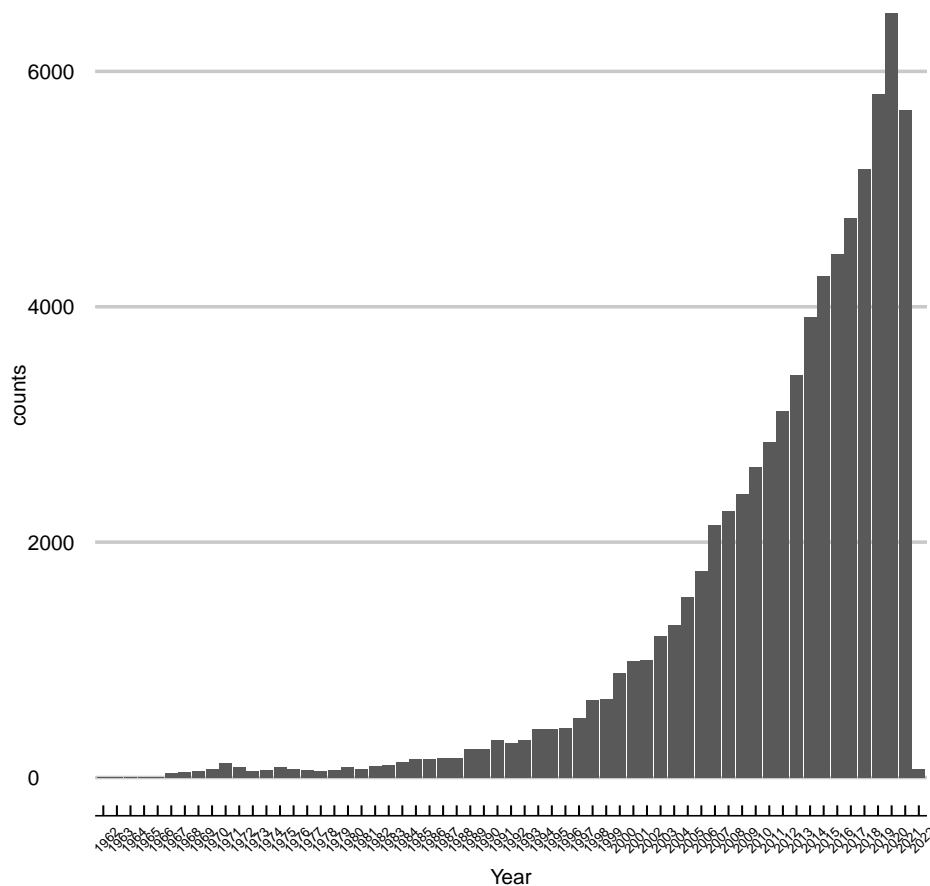


Figure 1: The searching result from PudMed.gov with key words :
Microbial Modeling

33 There are currently many mathematical models with different bene-
34 fits and drawbacks to illustrate the population growth curve[4]. The most
35 widely used mathematical models are the Logistic model[5],modified Gom-
36 pertz model[5] and Baranyi model[4]. Theoretically, if the per-capita growth
37 rate of a population is held constant (with limitless resource), the microbe
38 will lead to exponential unbound growth. However, it is not biologically

39 realistic. Therefore, one common way to address this deficiency is to use the
40 logistics model. In which, after the exponential growth, the rate will decrease
41 to zero as the population approaches a fixed value, also known as carrying
42 capacity. Although the logistic model is apparently simple and suitable for
43 some situations, it is still not generic enough in capturing other phenom-
44 ena. Therefore, the modified Gompertz model (Gompertz) [5] and Baranyi
45 model [4] which considered the influence of lag phase duration were created.
46 Although some literature had examined the goodness of fit among several
47 models [6, 7], most of them focus on the growth curve in one microbial species
48 which might be not generic for other species or under other conditions.

49 Therefore, Which model's performance is better to describe the general
50 microbial population growth is the main concern of this study. Based on 285
51 different empirical experiments data collected around the world with differ-
52 ent traits like species, growing medium, growing temperature etc., this study
53 focused on fitting and compare 6 different models (including linear models
54 and Non-linear models) to all the lab experiments data. Both Akaike infor-
55 mation criterion (AIC) and Bayesian information criterion (BIC) [8] were
56 used as criterion to assess the model fitting and to evaluate the performance
57 of model. However, given that the AIC and BIC performed 98% similar se-
58 lection, for the concern of efficiency and universality, AIC will be used as the
59 main criterion in this study. Moreover, Whether the environment (temper-
60 ature and medium), measurements (Units) will influence the performance
61 of the model will be assessed as an extra support to the model selection
62 strategy.

63 2 Data and Methods

64 2.1 Data Set and Preparation

65 The data set used in the study is called LogisticGrowthData.csv. The field
66 names are defined in a meta-file called LogisticGrowthMetaData.csv. Both
67 files are accessible in :

68 <https://github.com/nedchen2/CMEECourseWork/tree/master/MiniProject/data>.

69 The two main fields of interest are PopBio (abundance with different
70 units) and Time (Hours). Based on the combination identifier (unique
71 temperature-species-medium-citation combinations), 285 independent em-
72 pirical experiments were labelled with unique id from 1 to 285. Given the
73 log(exponential) phase during the growth, we did \ln transform to the abun-
74 dance (PopBio) and named the result as LogPopBio. Among the data sets,
75 there were some negative abundance which were replaced with a small num-
76 ber 0.00000000000000000001.

77 2.2 Mathematical Models

This study covers 6 models. In following equations, T is the time points
variables. $LogPopBio(T)$ is defined as the \ln (abundance) at the time point
 T

Simple Linear Regression model(named as Straightline in the report)

$$LogPopBio(T) = A_0 + A_1T + A_2T^2 \quad (1)$$

Quadratic polynomial model

$$LogPopBio(T) = A_0 + A_1T + A_2T^2 \quad (2)$$

Cubic polynomial model

$$LogPopBio(T) = A_0 + A_1T + A_2T^2 + A_3T^3 \quad (3)$$

Logistic model

$$LogPopBio(T) = \frac{N_0Ke^{r_{max}t}}{K + N_0(e^{r_{max}t} - 1)} \quad (4)$$

Gompertz model

$$LogPopBio(T) = N_0 + (K - N_0)e^{-e^{\frac{r_{max}e^1(t_{lag}-T)}{(K-N_0)\log(10)}+1}} \quad (5)$$

Baranyi model[4]:

$$LogPopBio(T) = N_0 + r_{max}(T + \frac{1}{r_{max}}\log(e^{-r_{max}T}) + e^{-r_{max}t_{lag}} - e^{-r_{max}(T+t_{lag})}) - \frac{\log(1 + ((e^{r_{max}(T+\frac{1}{r_{max}}\log(e^{-r_{max}T}))} - e^{-r_{max}t_{lag}})e^{-r_{max}(T+t_{lag})}))}{r_{max}} \quad (6)$$

78 The parameters used in non-linear least squares model have biological
79 meanings. N_0 is the initial population size, r_{max} is the maximum growth
80 rate , and K is carrying capacity (maximum possible abundance of the
81 population), t_{lag} is the duration of the delay before the population starts
82 growing exponentially. The algorithm I used to estimate the parameters
83 above will be introduced in subsequent part.

84 2.3 Model Fitting

85 With the development of computer and model fitting, multiple software can
 86 be applied least square algorithm for calculating the parameters of different
 87 models including linear (OLS) and non-linear model (NLLS). Levenberg-
 88 Marquardt algorithm [9] was applied to search and optimal parameter es-
 89 timates (or minimize the residual sum of squares,i.e. RSS). Maximum
 90 iteration number was set to be 200.

91 The start value estimate algorithm required in the non-linear model will
 92 be listed in table 2:

Table 1: The algorithms to estimate the start value

Parameters	Algorithm
N_0	minimum population of certain experiment
K	maximum population of certain experiment
r_{max}	the slope of straight line model fitting of certain experiment
t_{lag}	the time where the maximum differentials of population takes place

93 2.4 Model Selection

94 For the sake of the high accuracy, adjusted R-squared would be applied to
 95 compare the Linear Models[10]. However, when compares the non-linear
 96 and linear models, AIC and AIC_c and BIC of all models were calculated
 97 using following equations[8, 11]:

Akaike information criterion (AIC)

$$AIC = -2\ln(Likelihood) + 2k \quad (7)$$

Bayesian information criterion (BIC)

$$BIC = -2\ln(Likelihood) + k\ln(n) \quad (8)$$

$$AIC_c = AIC + \frac{2k(k+1)}{n-k-1} \quad (9)$$

Where n is sample size, k is number of parameters in the model, and $Likelihood$ is maximized likelihood function value of the model. AIC and BIC take both goodness of fit and model complexity into consideration to evaluate the performance of the model[8]. Although the formula of AIC and BIC are quite similar, they are distinguishing in theoretical bases[8]. The criterion was compared based on the empirical data considering the universality and efficiency to select the most appropriate one.

Model selection strategy.

Assuming that we are comparing a set of models with AIC values, we need a new data called ΔAIC which equals to the difference between a given model and the model with the lowest AIC.

Rough AIC: When one model's AIC equals to the minimum of the whole sets of AIC in different models, these models are identified as the best models.

Strict AIC: When one model's ΔAIC is less than 2[12], these models are identified as the plausible best models.

2.5 Extra Factors to Support The Model Selection

Measurements: In the data set, there are 4 different units which means the experiments were recorded by different measurements. Therefore, we want

117 to know that if the measurement is relevant to the model performance.

118 **Temperature:**In original data set, the temperatures ranged from 0 to 37.

119 In order to better categorize the temperature[13], we roughly classify the

120 temperature between 0 - 5 into "cold lover", the temperature between 5-20

121 into "middle cold lover", and the temperature between 20-37 into "middle

122 hot lover".

123 **Medium:**Finally, we classify the medium which is more routine such as milk

124 and chicken as "Nature", and the medium which is much more lab-based

125 like TSB and MRS as "Artificial".

126 After classification, the best fitted number of six types of models in ev-

127 ery category was counted. Then the frequency would be transformed into

128 percent and plotted to bar plot.

129 2.6 Computing Software

130 For the reason of powerful and convenient application, R 4.1.1 were used as

131 the main computing tools for data wrangling,model fitting and visualization.

132 Several packages such as ggthemes, ggpubr, tidyverse (for data wrangling

133 and visualization), and minpack.lm (for model fitting) were imported in

134 the study for different purpose. A bash script was used to compile LaTeX.

135 Python3, particularly with the subprocess package, was used as a work flow

136 control tool to run all the scripts.

3 Results

3.1 Linear Least Squares Model Fitting

Simple linear regression model (straight line), quadratic polynomial, cubic polynomial were applied to fit the data. Adjusted R-squared are used as the criterion to evaluate the model fitting.

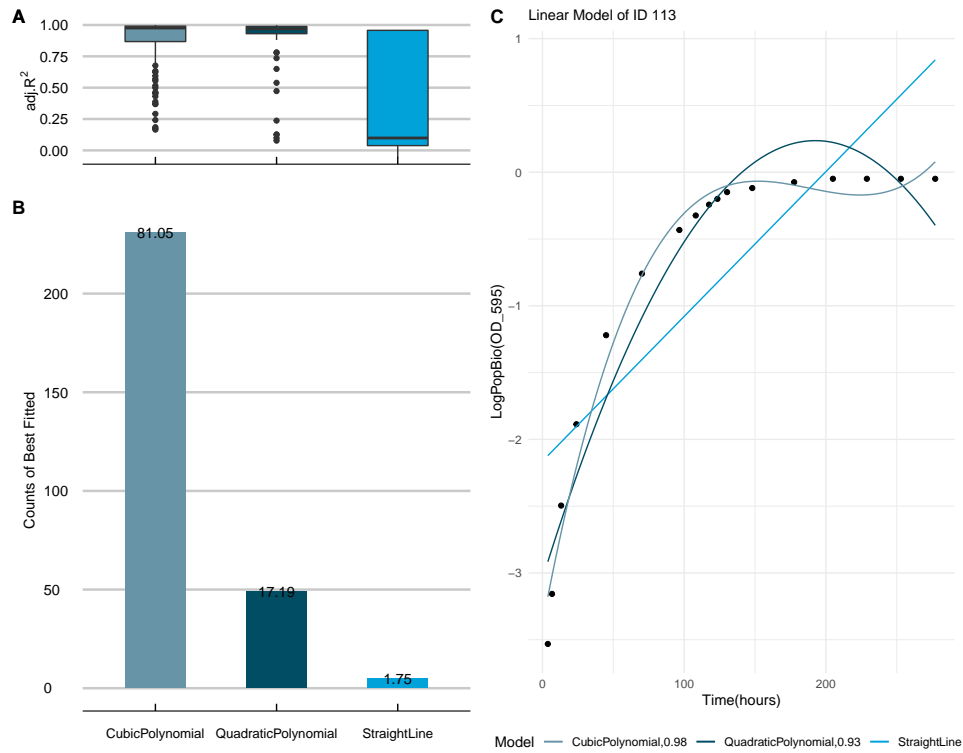


Figure 2: Linear model fitting. A. The box plot to show the range of adjusted R-squared in different models. B. The bar plot to show the best fitted counts in different models. C. An example of linear fitting

In Figure 2.A, the diagram shows that the cubic polynomial and quadratic polynomial model had a higher fitted R-squared value (close to 1). In the Figure 2.B, the bar plot shows the distribution of 3 different models fre-

quency which best fits the 285 experiments. As demonstrated in the figure, the counts of best fitted cubic polynomial is much larger than other two models, which accounts for 81.05% percent of the 285 experiments. In the Figure 2.C, the figure shows the fitting of three linear models of experiment ID 113. The adjusted R-squared value are displayed after the model name. The cubic polynomial fitting shows the best fitting, while the Straightline Model shows the worst fitting. However, Figure 2.C also illustrates that the cubic polynomial model can not fit the stationary phase properly, despite that the R-squared are 0.98.

3.2 Comparison between AIC, AIC_C , BIC

In order to compare the OLS with NLLS model, we will use AIC, AIC_C or BIC as criterion in the comparison and selection. To simplify the analysis procedure, I decided to compare the three criterion first. A criterion which could best represented the model selection result would be chosen as the criterion for subsequent research.

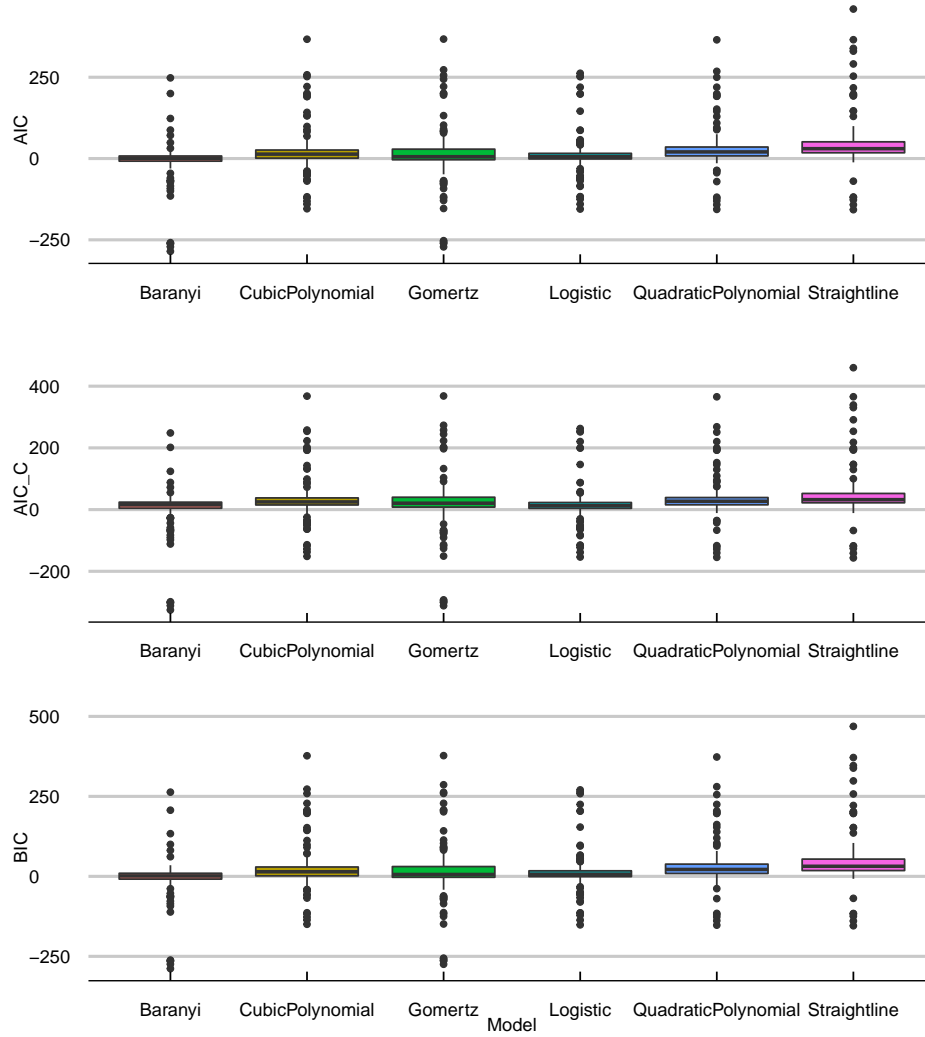


Figure 3: Box plot of different models show the range of AIC,BIC,AIC_C.
The colours of Box represent different plausible models

160 Generally, in all criteria, Straightline model shows the largest average
 161 value. NLLS models including logistics,Gompertz,Baranyi models shows
 162 lower average value. Consequently, the figure in different criteria shows the
 163 similar pattern of value distribution. After Comparing the model selection

164 result based on rough AIC algorithm, the result shows that 97.54% of the
 165 best fitted result are supported both by AIC and BIC. However, only about
 166 50% of the best fitted result are supported both by AIC and AIC_C . For the
 167 universality and efficiency of the study, the decision was made to use AIC
 168 in the subsequent report.

169 3.3 Linear and Non-linear Model Fitting

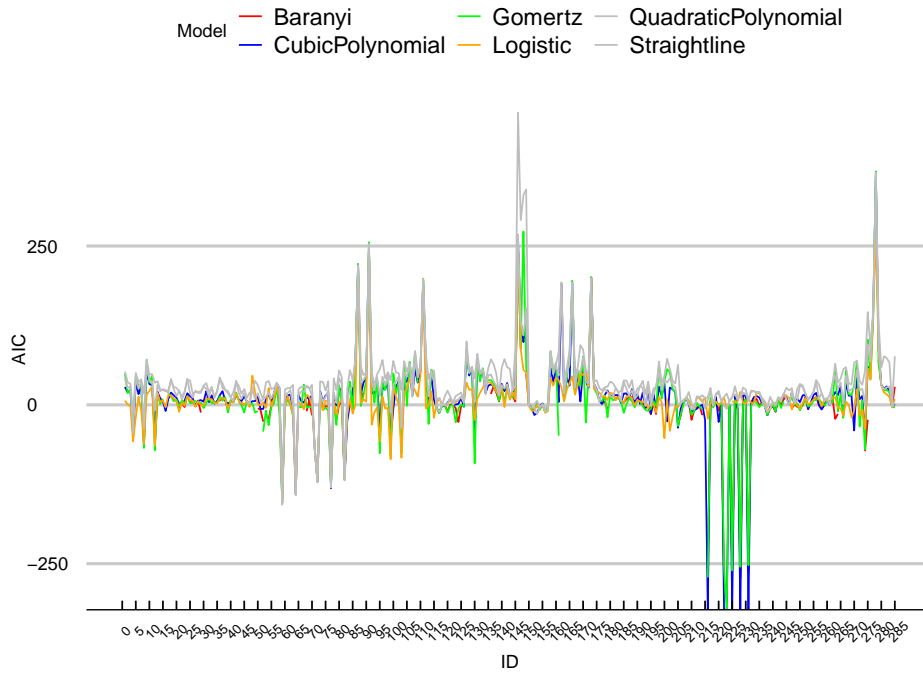


Figure 4: AIC distribution of different models across the experiment ids.
 The colours of lines represent different plausible models

170 Figure 4 shows the AIC value distribution of different models across the
 171 experiment ids. Based on the result of the preceding linear model fitting,
 172 the misfitted quadratic polynomial model and the straight line model will

be defined as grey lines in Figure 4. As illustrated in the figure, a large portion of the yellow curve shows the lowest AIC value across the IDs axis. However, implicit difference in distribution of the AIC value in the various models can also be observed.

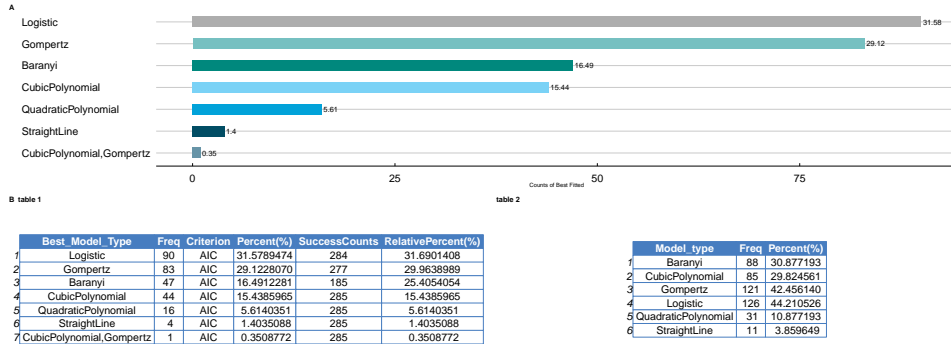


Figure 5: Best Fitted Models. A. The bar plot of best fitted models based on rough AIC criteria. The number illustrated beside the bar is percentage (counts of best fitted/total experiments number); B. table 1 is the table based on rough AIC. Freq column shows the counts of best fitted model. SuccessCounts column shows the successfully modelling counts for each models. RelativePercent column shows the relative best fitted percentage (Freq/SuccessCounts); B. table 2 is the table based on strict AIC

During the modelling, some data set could not be fitted by NLLS model. The reason might be that the start value is too far from the correct answer. Therefore, relative best fitted percent which indicates the best fitted counts over the successfully modelled counts will be used to generalize this issue, as illustrated in Figure 5.B(table1). Figure 5(A) shows that based on rough AIC selection algorithm, 3 non-linear models (Logistic model, Gompertz model, Baranyi model) are the top 3 best fitted models. In one experiments, both cubic polynomial and Gompertz are considered as the best model. 31.93% of the experiments fits the Logistic model best. 29.12% of the experiments fits the Gompertz model best. 16.14 % of the exper-

187 iments fits the Baranyi models best. 15.44% of the experiments fits the
 188 cubic polynomial best. Although the difference between the performance of
 189 cubic polynomial and Baranyi model is small, considering the relative best
 190 fitted percentage (demonstrated in Figure 5.B(table 1), the performance of
 191 Baranyi model is much better than cubic polynomial. The Figure 5.B(table
 192 2) shows the result based on strict selection algorithm. The result shows
 193 that the best fitted percent of Logistic and Gompertz model raises dramati-
 194 cally, indicating that the performance of these two models are much better
 195 than other models.

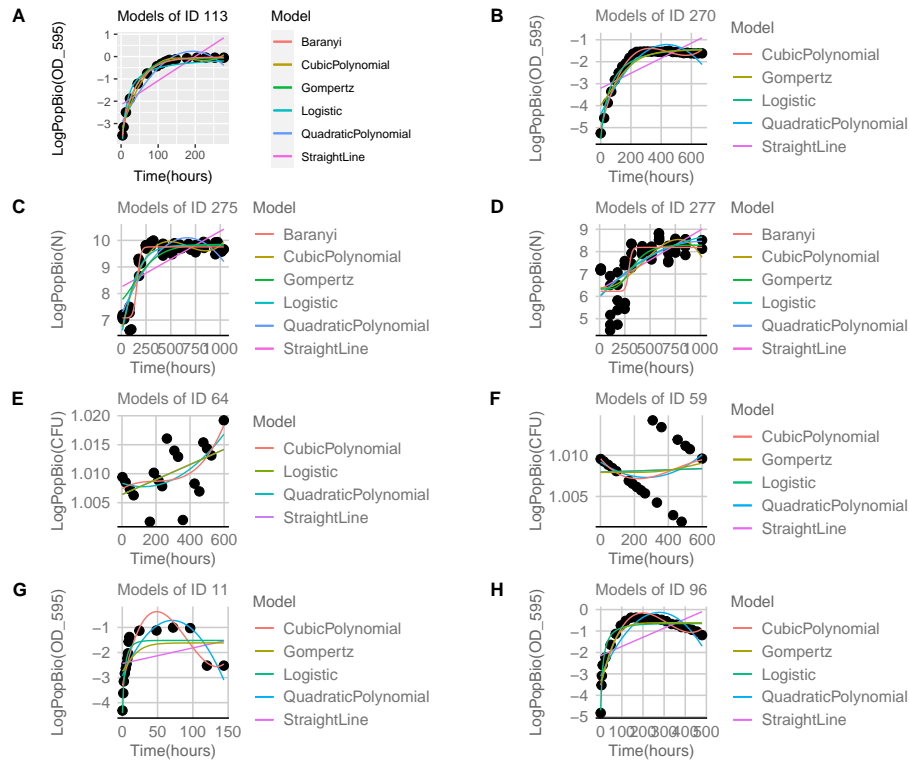


Figure 6: Examples of Model fitting in 285 experiments. Different models are represented in different colours.

196 The Figure 6 A,B show logistic growth phase curve and stationary phase

197 curve. The difference of various NLLS models is not explicit in these cases. In
198 the Figure6 C,D, sigmoid growth curve is displayed. We can also find that
199 the non-linear model including Baranyi, Gompertz models performs much
200 better when the lag phase exists. However, although the Baranyi and Gompertz
201 model can fit the lag phase, exponential phase and stationary phase,
202 they can not fit the death phase properly, as illustrated in Figure 6 G,H. In
203 the Figure6 E,F, the scatter shows unexpected pattern. In both cases, the
204 straightline model shows the lowest AIC.

205 **3.4 Models Performance of Different Measurements, Temperature, Medium**

206

207 After classifying the best fitted result into different categories, the percentage
208 of each models in different category was displayed.

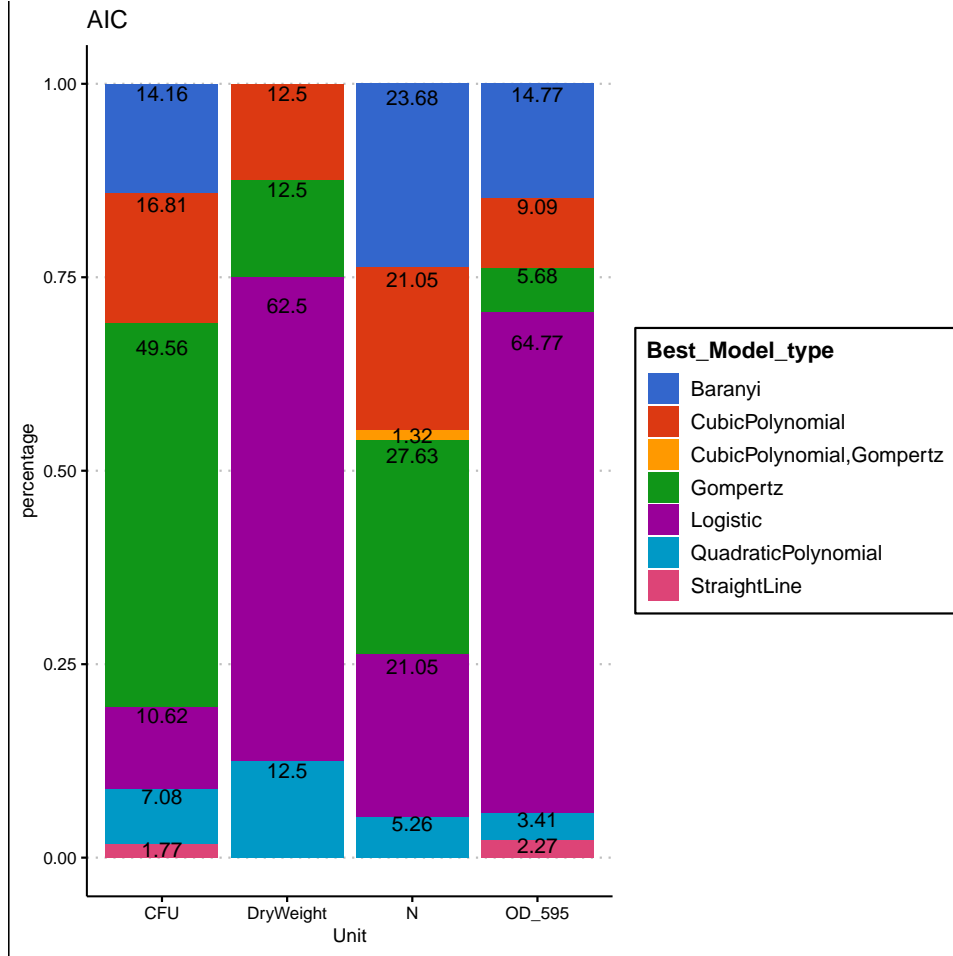


Figure 7: Best Fitted Model in different measurements. Different models are represented in different colours

Figure 7 shows that logistic models perform well (accounting for over 60% of best fitted model) for OD595 and DryWeight, while for CFU, Gompertz performs well.

According to previous work[?], we roughly classify the temperature between 0 - 5 into "cold lover", the temperature between 5-20 into "middle cold lover", and the temperature between 20-37 into "middle hot lover".

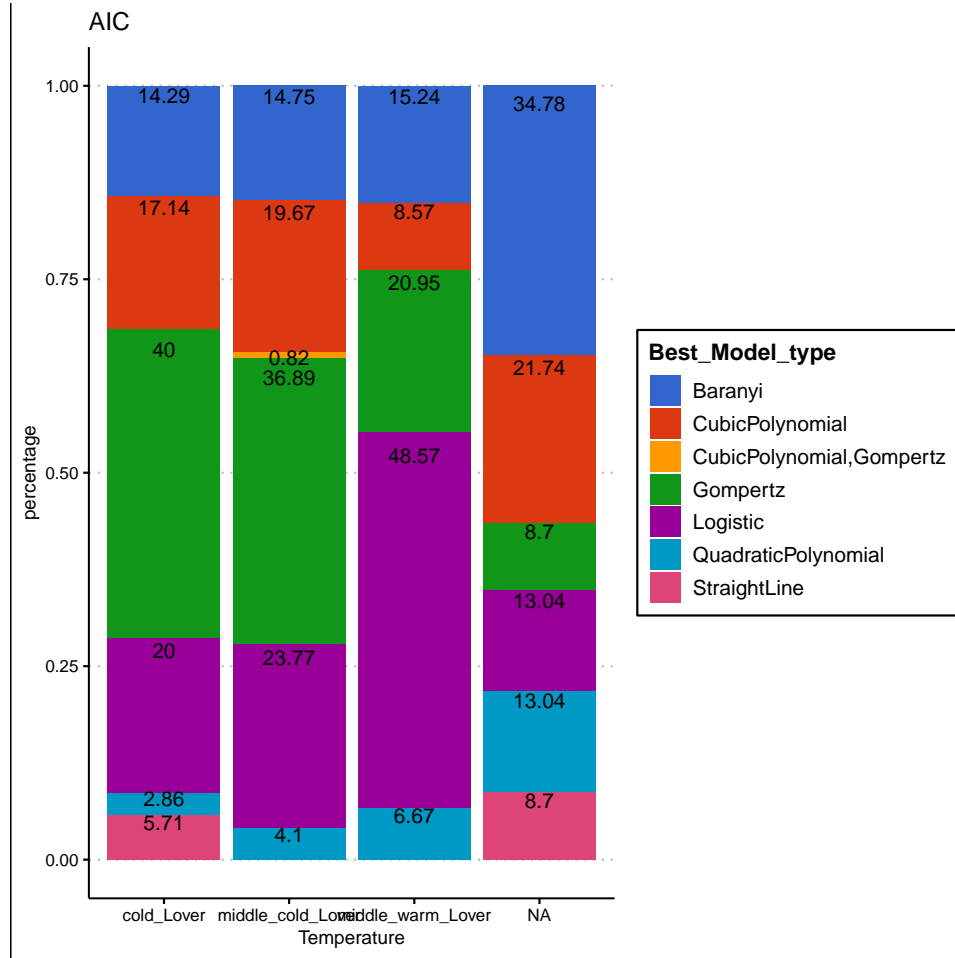


Figure 8: Best Fitted Model in different Temperatures. Different models are represented in different colours

215 Generally, Logistic model performs better in middle warm. Moreover,
 216 as the temperature grows, the best fitted percentage accounted by logistic
 217 model tends to grow. Additionally, the straightline model only displays in
 218 cold temperature.

219 Finally, we classify the medium which is more natural such as milk and
 220 chicken as "Nature", and the medium which is much more artificial like

221 TSB and MRS as "Artificial". After that, we count the best fitted models
 222 in "Nature" and "Artificial".

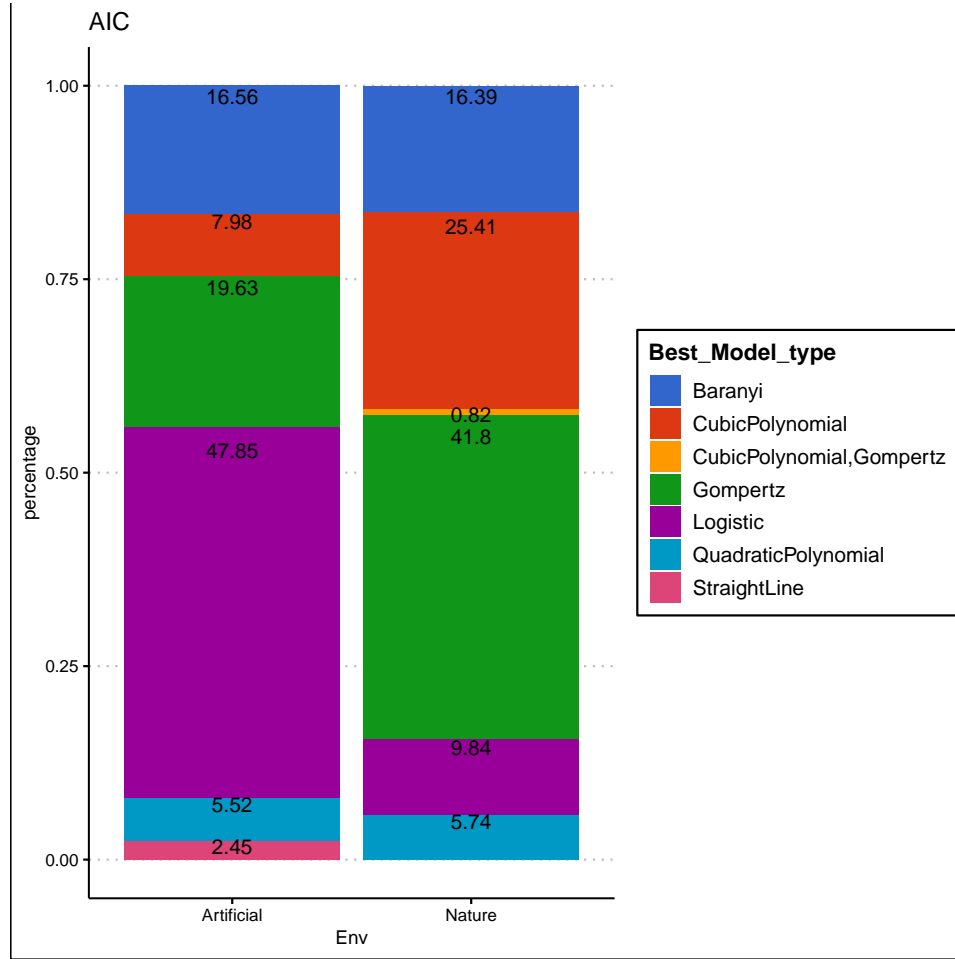


Figure 9: Best Fitted Model in different Medium. Different models are represented in different colours

223 The results shows that the Logistics model will be preferred in "Arti-
 224 ficial" medium, while the Gompertz models will be preferred in "Nature"
 225 medium.

226 4 Discussion

227 Understanding how microbial organisms grow under particular conditions is
228 essential for human to study population dynamics and microbe ecosystem[3].
229 In terms of microbial growth prediction and evaluation, choosing an univer-
230 sal and appropriate mathematical model which is supported by empirical
231 experiments and valid data is an inevitable issue[14]. Currently, many mod-
232 els are pointed out to fit the microbial growth model. Therefore, to test the
233 universality and application of certain model on empirical data, I did model
234 fitting and selection on 285 worldwide experiments data set by 6 typical
235 models.

236 NLLS Model fitting is challenging because of the start value choice. For
237 example, the r_{max} , I have tried 3 methods including using the maximum
238 derivatives of cubic polynomial model, using the A2 parameter of quadratic
239 polynomial model and the and using the slope of the simple linear regression
240 model. I choose the last method which have the highest successful fitting
241 rate through the 285 experiments data set. However, the Baranyi model's
242 successful fitting rate is still low (65%). Additionally, I did not randomly
243 sample the start value to optimize the result of fitting. Therefore, efforts
244 should be taken in the future to optimize the fitting by either improving the
245 algorithm of the start value, or random sampling the start value.

246 Based on the comparison of various model criterion, in our study, almost
247 98% of the best fitted result is both supported by AIC and BIC as the best
248 model. Therefore, for the concern of universality and efficiency, I decided to
249 use AIC as the main criterion in further research.

250 During model selection, I define two selection rules : rough AIC and
251 strict AIC to evaluate the performance of models. Firstly, regardless of the

rules, non-linear least squares models (include Logistic, Gompertz, Baranyi models) perform better than linear models. In terms of the linear model, cubic polynomial model performs the best.

The results also showed that under strict AIC selection rules, logistic model is sufficient in fitting microbial growth for about 44% of experiments in the data set. Gompertz model is sufficient in fitting microbial growth for about 40% of experiments in the data set. Given the low successfully modelled rate, Baranyi Models can also perform better if we could improve the model fitting by some methods such as changing the start value algorithms.

Investigating the result with the ideal microbial growth curve in a "closed habitat", we can find some drawbacks among all models, regardless of linear or non-linear model. For instance, although cubic polynomial is good to fit the log(exponential) phase in empirical data, the cubic polynomial curve can not fit the lag phase and stationary phase properly, as illustrated in Figure 6 C,D and Figure 2 C. Similarly, although the Baranyi and Gompertz model can fit the lag phase, exponential phase and stationary phase, they can not fit the death phase, as illustrated in Figure 6 G,H. Furthermore, Straightline model might be a good indicator of abnormal data set, as illustrated in Figure 6 E,F. Consequently, although no models are universal for all conditions, given the simple and efficient formula, logistic models can be considered as the sufficient model for most cases. However, logistic models can not handle the lag phase and death phase. Therefore, combine the model selection strategy with segmented growth phase might be a solution for handling different situations.

Apart from the segmented model strategy, we may use other factors to assist the model selection. For example, as illustrated in Figure 7, the logistic models are preferred (65%) when using OD595 as the microbial pop-

ulation measurements. However, the reason is unknown. Another example is that, the logistic models are preferred when the temperature getting warmer, which may be partly explained by that when the temperature increasing, the growth pattern may be more similar to logistic growth[13]. Finally, we classify the medium which is more natural such as milk and chicken as "Nature", and the medium which is much more artificial like TSB and MRS as "Artificial". The results shows that the Logistic model will be preferred in "Artificial" medium, while the Gompertz models will be preferred in "Nature" medium. The phenomenon might be explained by that in the "Nature" medium, the lag phase would be more apparent which is preferred by the Gompertz models.

In conclusion, based on the model fitting and model selection on empirical data sets, non-linear model performs better than linear model. Among the non-linear model, logistic model is sufficient for simple growth situation, while Gompertz and Baranyi can handle more complicated situation like lag phase. However, all of the non-linear models such as logistic model and Baranyi model have the defect that they can not be applied to fit the death phase properly. Segmented model and external factors calibration might be a potential strategy to optimize the result.

5 Supplementary Materials

Please find the complete model fitting plot of every experiments in

<https://github.com/nedchen2/CMEECourseWork/tree/master/MiniProject/results/>

Filename Format: ID + ALL Model Plot.png

Please find the model fitting result data for each models in

<https://github.com/nedchen2/CMEECourseWork/tree/master/MiniProject/results>

304 Filename Format: Model type + result.csv
305 Please find the complete best model result depending on AIC,BIC, AIC_C in
306 <https://github.com/nedchen2/CMEECourseWork/tree/master/MiniProject/results>
307 Filename Format: Best Model + criteria.csv
308 Please find the code in
309 <https://github.com/nedchen2/CMEECourseWork/tree/master/MiniProject/code>

310 References

- 311 [1] William R. Shoemaker, Stuart E. Jones, Mario E. Muscarella, Megan G.
 312 Behringer, Brent K. Lehmkuhl, and Jay T. Lennon. Microbial popu-
 313 lation dynamics and evolutionary outcomes under extreme energy lim-
 314 itation. 118(33). Publisher: National Academy of Sciences Section:
 315 Biological Sciences.
- 316 [2] Thomas Egli. Microbial growth and physiology: a call for better crafts-
 317 manship. 6:287.
- 318 [3] Micha Peleg and Maria G. Corradini. Microbial growth
 319 curves: What the models tell us and what they can-
 320 not. 51(10):917–945. Publisher: Taylor & Francis .eprint:
 321 <https://doi.org/10.1080/10408398.2011.570463>.
- 322 [4] J. Baranyi, P. J. McClure, J. P. Sutherland, and T. A. Roberts. Mod-
 323 eling bacterial growth responses. 12(3):190–194.
- 324 [5] M. H. Zwietering, I. Jongenburger, F. M. Rombouts, and K. van ’t Riet.
 325 Modeling of the bacterial growth curve. 56(6):1875–1881.
- 326 [6] R Xiong, G Xie, A. S Edmondson, R. H Linton, and M. A Sheard.
 327 Comparison of the baranyi model with the modified gompertz equation
 328 for modelling thermal inactivation of listeria monocytogenes scott a.
 329 16(3):269–279.
- 330 [7] R. L Buchanan, R. C Whiting, and W. C Damert. When is simple good
 331 enough: a comparison of the gompertz, baranyi, and three-phase linear
 332 models for fitting bacterial growth curves. 14(4):313–326.

- 333 [8] Jouni Kuha. AIC and BIC: Comparisons of assumptions and perfor-
334 mance. 33(2):188–229. Publisher: SAGE Publications Inc.
- 335 [9] Henri P Gavin. The levenberg-marquardt algorithm for nonlinear least
336 squares curve-fitting problems. page 19.
- 337 [10] Davide Chicco, Matthijs J. Warrens, and Giuseppe Jurman. The coef-
338 ficient of determination r-squared is more informative than SMAPE,
339 MAE, MAPE, MSE and RMSE in regression analysis evaluation.
340 7:e623.
- 341 [11] Eric-Jan Wagenmakers and Simon Farrell. AIC model selection using
342 akaike weights. 11(1):192–196.
- 343 [12] Adrian E. Raftery. Bayesian model selection in social research. 25:111–
344 163. Publisher: [American Sociological Association, Wiley, Sage Pub-
345 lications, Inc.].
- 346 [13] Andrea De Silvestri, Enrico Ferrari, Simone Gozzi, Francesca Marchi,
347 and Roberto Foschino. Determination of temperature dependent
348 growth parameters in psychrotrophic pathogen bacteria and tentative
349 use of mean kinetic temperature for the microbiological control of food.
350 9:3023.
- 351 [14] Kathleen M. C. Tjørve and Even Tjørve. The use of gompertz models
352 in growth analyses, and new gompertz-model approach: An addition to
353 the unified-richards family. 12(6):e0178691. Publisher: Public Library
354 of Science.