

High-Throughput Non-Linear Model Comparison and Analysis for Bacterial Growth Curves: a Bayesian Approach

Ben Nouhan, bjn20@ic.ac.uk

Imperial College London

January 22, 2021

Word Count: 3374

1 The ability to understand and predict population growth is vital for multiple disci-
2 plines. Technology is increasingly enabling us to model large datasets, uncover the
3 insights buried within them, and improve the models iteratively. Scaling up this
4 process to ingest more data and quantify improvements to the models would push
5 forward our capabilities and inform future research. Here I showcase a prototypal
6 pipeline to fit established models to hundreds of datasets, quantify their perfor-
7 mance for comparison between them, and use control variable data to glean insights
8 out of this process. Consistent with the literature, the methodology proclaimed the
9 Gompertz model as the highest-performing of those tested, while highlighting its
10 flaws. Correlating performance of the models and separately the morphology of
11 their resultant fits with potential covariables has the potential to improve or even
12 inspire subsequent investigation. Meanwhile the pipeline as a whole can, with
13 modest alterations, be used on groups of models from a multitude of fields, at best
14 facilitating the development of the very models it upon which it is used to analyse
15 and elucidate.

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 2 | Methods | 4 |
| 2.1 | Computing Tools | 4 |
| 2.1.1 | Python | 4 |
| 2.1.2 | R | 4 |
| 2.1.3 | Resources | 4 |
| 2.2 | Data | 5 |
| 2.2.1 | Raw Dataset | 5 |
| 2.2.2 | Preprocessing | 5 |
| 2.3 | Model Fitting | 5 |
| 2.3.1 | Linear Models | 5 |
| 2.3.2 | Non-Linear Models | 5 |
| 2.4 | Model Comparison | 6 |
| 3 | Results | 7 |
| 3.1 | Model Fitting | 7 |
| 3.2 | Model Comparison | 7 |
| 3.3 | Model-Covariable Correlation | 7 |
| 4 | Discussion | 11 |
| | References | 13 |

1 Introduction

Understanding population growth is paramount in fields of study as far-flung as epidemiology, climate science and geopolitics.(Ozgul et al., 2010; Peleg, 1997) For decades, increasingly complex mathematical models have been used to explain trends in empirical population growth time series and enable prediction.(Grijnspeerdt & Vanrolleghem, 1999; Kingsland, 1982; Tjørve & Tjørve, 2017) Fewer parameters reduce the chance of models overfitting the data and hence, using bacterial growth as an example, variables not included such as incubation temperature, bacterial strain and growth medium should be kept constant.

Bacterial growth models largely rely on the theory of bacterial growth phases in a closed system, shown schematically in **Figure 1**. There are four accepted phases: the lag phase, exponential growth phase, stationary phase and death phase, with some considering the three transition periods between them as additional phases in their own right.(R. E. Buchanan, 1918)

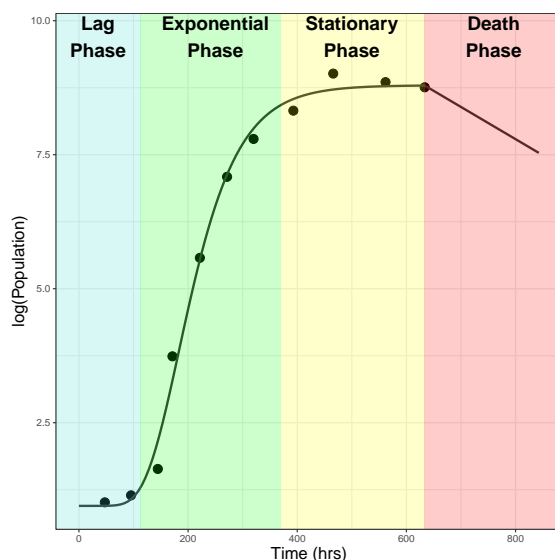


Figure 1: Schematic example of an archetypal bacterial growth curve demonstrating the four phases of growth. The lag phase is the initial period of zero or minimal growth whereby the bacteria, having been transferred to a new medium, require time to acclimatise. For example, the new environment may impact gene expression and hence the bacteria's replication machinery are not immediately operational.(R. E. Buchanan, 1918) During the exponential growth phase, characterised by an exponential curve owing to the rate of increase per bacterium remaining constant, the bacteria can continuously multiply absent limiting environmental factors.(Peleg & Corradini, 2011) The stationary phase, a population plateau completing the sigmoidal shape of the growth curve, arises due to the population reaching the carrying capacity of the medium; rate of division approximately equals the death rate.(R. E. Buchanan, 1918) In some instances there is a subsequent death phase, during which death rate surpasses the rate of division due to factors such as the accumulation of a toxic substance or depletion of the medium.(R. E. Buchanan, 1918; Peleg & Corradini, 2011) The first three phases in the above chart were generated from the dataset used throughout this study (specifically, *Pseudomonas* spp. grown on raw chicken breast at 2°C) with a regression line fit using the Gompertz model. The "death phase" was appended artificially.

Technological advancements since the 1990s have caused the quantity of data generated from biological experiments, the speed at which computers analyse them and the accessibility of the process to skyrocket. This allowed life scientists to mathematically model natural phenomena in ways previously limited to the physical sciences.(Bolker et al., 2013; Johnson & Omland, 2004)

These can be linear models (LMs), wherein the response variable 'y' has a linear relationship with the parameters if not the explanatory variable 'x' as in **Equation 1**; or non-linear models (NLMs), wherein the response variable has a non-linear relationship with a parameter and the explanatory variable as in **Equation 2**. (Bolker et al., 2013) Regardless of whether these models truly represent the natural laws in question, they are undoubtedly useful for prediction phenomena, and the development of more sophisticated models. (Transtrum & Qiu, 2016)

$$y = a + bx + cx^2 \quad (1)$$

$$y = a + bx^c \quad (2)$$

Parameters of some NLMs for population growth, such as the logistic, Gompertz, Baranyi and Buchanan models, can be related to the aforementioned phases. These include: t_{lag} , the duration of the lag phase; N_0 , the minimum population that can feasibly lead to growth; N_{max} , the maximum population the system can feasibly support; and r_{max} , the maximum possible growth rate. (Peleg & Corradini, 2011) It has been asserted that without parameters like these based firmly in scientific theory an equation used to fit data is not truly a model. (R. L. Buchanan et al., 1997)

The logistic model, **Equation 3**, is one of the oldest population growth models and is still used in fields from medicine to economics. It was initially posited as a model for human population growth in which the growth rate per unit decreases as the sample population approaches N_{max} . (Peleg, 1997) Many newer, more sophisticated population growth models were derived from the logistic model, but introduce the t_{lag} parameter which increases their utility when fit to timeseries with a lag phase.

$$N_t = \frac{N_0 \cdot N_{max} \cdot e^{t \cdot r_{max}}}{N_{max} + N_0 \cdot (e^{t \cdot r_{max}} - 1)} \quad (3)$$

The modified Gompertz model incorporates biologically meaningful parameters into an empirical, sigmoidal relationship. (R. L. Buchanan et al., 1997; Tjørve & Tjørve, 2017) First conceived for predicting mortality rates in human populations, countless studies in various disciplines have utilised it. (R. L. Buchanan et al., 1997; Mokhtari et al., 2019; Peleg, 1997; Tjørve & Tjørve, 2017) One form of it, shown by **Equation 4**, includes the t_{lag} parameter, thereby incorporating the lag phase absent in the logistic model. (R. L. Buchanan et al., 1997; Zwietering et al., 1990)

$$N_t = N_{max} \cdot e^{-e^{\frac{e \cdot r_{max}}{N_{max} - N_0} \cdot (t_{lag} - t) + 1}} \quad (4)$$

The Baranyi model, first published in 1993, is a logistic rate differential equation designed

specifically for modelling bacterial growth curve dynamics.(Baranyi et al., 1993; R. L. Buchanan et al., 1997) The theory of a "bottleneck" chemical reaction limiting the maximum growth rate, r_{\max} , underpins the model.(R. L. Buchanan et al., 1997) Alongside the Gompertz model it has overtaken the logistic model in popularity for modelling population growth, owing in part to the t_{lag} parameter that can be derived from the original equation: **Equation 5** represents the baranyi model rearranged to include the parameters discussed herein.

$$N_t = N_{\max} - \ln(1 + (e^{-N_{\max} - N_0} - 1).e^{-r_{\max}.t_{\text{lag}}}) \quad (5)$$

The Buchanan model can be thought of as a three-phase linear model, demonstrated by **Equation 6**.(R. L. Buchanan et al., 1997) It was proposed in 1997 to determine how accurately bacterial growth timeseries could be modelled by a simpler model to those of Gompertz and Baranyi. It requires a parameter t_{\max} , the time at which N_{\max} is first reached, which can be estimated from N_{\max} itself. Its first phase exhibits zero growth until approximately t_{lag} , preceding a period of linear r_{\max} growth, until the population plateaus at t_{\max} .(R. L. Buchanan et al., 1997)

$$N_t = \begin{cases} N_0 & \text{for } t \leq t_{\text{lag}} \\ N_{\max} + r_{\max} * (t - t_{\text{lag}}) & \text{for } t_{\text{lag}} < t < t_{\max} \\ N_{\max} & \text{for } t \geq t_{\max} \end{cases} \quad (6)$$

Modelling is seen by some as the successor of classical hypothesis testing, and by others as another tool in their arsenal.(Johnson & Omland, 2004) Since multiple models can be employed for the same task, the ability to determine which model is most useful in a given situation is a science in of itself. The Bayesian information criterion (BIC) is one metric for model selection, which is defined in **Equation 7**. It assigns each model a score derived from the sample size used, n , the maximum likelihood estimation of the model, L , and the number of parameters, k ; models with fewer parameters are generally more stable by reducing inter-parameter correlation.(Akaike, 1974; Zwietering et al., 1990) The lower the score, the better the model; a difference in score between models of less than two is considered insignificant, and greater than ten highly significant.(Posada & Buckley, 2004; Vrieze, 2012). An alternative is Akaike information criterion (AIC) which, while derived from frequentist probability rather than Bayesian, is largely the same except confers a smaller penalty for additional parameters.(Posada & Buckley, 2004)

$$BIC = k.\log(n) - 2.\log(L) \quad (7)$$

While BIC scoring can inform on the relative performance of multiple models on a given dataset, there is no established method for doing this across multiple datasets. Thus my objectives with

89 this study are three-fold: to design a general, robust methodology for the high-throughput fitting
90 of multiple population growth models, linear and non-linear, to a large quantity of datasets; to
91 further design a method for selecting the overall best model, determined as a function of both
92 accuracy and consistency; and to visualise the results in a way that will highlight correlations
93 between covariables of the datasets and performance of the models or the growth behaviour. The
94 latter may reveal whether certain models may be more appropriate for experiments executed under
95 certain conditions, or if said conditions alter the underlying mechanism being modelled, facilitating
96 the conception of hypotheses for subsequent experimentation.

97 **2 Methods**

98 **2.1 Computing Tools**

99 **2.1.1 Python**

100 Data preprocessing was performed using Python 3.9.0 for its superior data processing functionality.
101 Its *pandas* package enables efficient and user-friendly database manipulation, while its *numpy* package
102 enables generation of unique hash identifiers and log transformation of data.

103 **2.1.2 R**

104 R 4.0.3 was utilised for model fitting and subsequent analysis and visualisation of results. Ar-
105 guably less generally capable than Python, it was built specifically for statistical analysis, hence
106 the tools required are presently more established, comprehensive and supported in R than their
107 Python counterparts. Additionally, the R Core distribution includes the *parallel* package allowing
108 the harnessing of all available computer cores for more computer-intensive tasks which, alongside
109 vectorisation of model fitting and plotting, cuts the runtime to mere seconds. The *tidyverse* package
110 is needed for a wealth of processes, from efficiently importing dataframes to increasing dataframe
111 manipulation and data visualisation capabilities.

112 **2.1.3 Resources**

113 The dataset used to develop the workflow, alongside explanatory metadata and the workflow itself,
114 are available on github ([github.com/mhasoba/TheMulQuaBio/blob/master/content/data/L](https://github.com/mhasoba/TheMulQuaBio/blob/master/content/data/LogisticGrowthData.csv)
115 [ogisticGrowthData.csv](https://github.com/mhasoba/TheMulQuaBio/blob/master/content/data/LogisticGrowthData.csv), [github.com/mhasoba/TheMulQuaBio/blob/master/content/data](https://github.com/mhasoba/TheMulQuaBio/blob/master/content/data/LogisticGrowthMetaData.csv)
116 [/LogisticGrowthMetaData.csv](https://github.com/Bennouhan/cmeeecoursework/tree/master/minipr) and github.com/Bennouhan/cmeeecoursework/tree/master/minipr
117 object respectively). See *README.md* for details and dependencies.

2.2 Data

2.2.1 Raw Dataset

The dataset comprises 305 bacterial population growth timeseries from a multitude of published studies in a long-format, 4388 row dataframe, each row representing a datapoint. These timeseries use a variety of variables, each constant for a given timeseries: 17 incubation temperatures irregularly spaced between 0 and 37°C; 18 growth media; 45 bacterial species; and four population estimation techniques. These were colony-forming units (CFU) counts, a sample's dryweight, sample optical density at OD-595, and CFUs of differently appearing colonies in a mixed-species sample (referred to as 'N' in the dataset). (Al-qadiri et al., 2008) For later visualisation nine growth media, used in fewer than ten timeseries, were excluded.

2.2.2 Preprocessing

The workflow organises and cleanses the dataset before indexing each timeseries with a unique ID, facilitating subsequent referencing. This includes calibrating timeseries containing a negative initial time measurement to zero since these are likely systematic errors, and deleting negative population datapoints, assuming them irreconcilable errors. The population measurements are subsequently $\log_2(x+1)$ transformed. Taking the log of the population data makes processing and visualising their wide range between studies more intuitive by normalising them and their measurement unit, and is to no detriment since only relative changes in population are of interest. Base 2 is used because bacteria duplicate, while the $\log(x+1)$ transformation accepts population measurements below 1. Fifteen timeseries with fewer than six datapoints were omitted to avoid the LMs with up to five parameters overfitting those timeseries.

2.3 Model Fitting

2.3.1 Linear Models

I performed linear regression by fitting each of the time series with four LMs of first, second, third and fourth order polynomials, using R's *lm()* function. These correspond to linear, quadratic, cubic and quartic expressions respectfully, named as such in this paper. **Equation 1** is the generic form of a second order, quadratic model, with three parameters, a-c. The Linear, cubic and quartic expressions have two, four and five parameters respectively, following the sequence. These parameters are found by ordinary least square (OLS) linear regression and not based upon scientific theory, hence LMs are phenomenological unlike mechanistic NLMs.

2.3.2 Non-Linear Models

I fitted the NLMs to the time series using the *nlsLM()* function from the *minpack.lm* R package. With the exception of the logistic model, lacking t_{lag} , all four of the NLMs used - logistic, Gompertz,

Baranyi and Buchanan - have identical parameters underpinned by the same scientific theory and can, therefore, be interpreted uniformly.(Odenbaugh, 2006)

As mechanistic models, NLM fitting requires parameter starting values. If sufficiently close to the true value, the function will perform fitting iteratively until converging on the function's parameter estimates. To adequately predict the true parameter values for each timeseries, I utilised their datapoints. Within each timeseries, my estimate for: N_{\max} was the highest population present; N_0 was the lowest population present; r_{\max} was the highest rate of the change between two adjacent datapoints; and t_{lag} was that same rate of change extrapolated from those datapoints to the intercept with my N_0 estimate.(Peleg & Corradini, 2011)

R^2 , a common measure of correlation, is invalid for NLMs as the sum of squares of residuals divided by the total sum of squares can exceed 1, nor does it consider model complexity thereby ignoring the principle of parsimony.(Johnson & Omland, 2004) However, it can be used as a rough proxy of correlation or to compare multiple fits solely in that regard, hence I used the R^2 of the resulting regression line of each timeseries to assess the fit. If the fit had failed or had an R^2 below 0.5, I reattempted the fit up to 500 times by randomly selecting starting values from a normal distribution around the original estimate with a standard deviation of four times that original estimate.

2.4 Model Comparison

I calculated the BIC score and R^2 of each successful fit before using them to compare models by their collective fits. To remove poor datasets, those where no model's fit had an R^2 above 0.7 were excluded, leaving 256. Since the NLMs, excluding logistic, have equinumerous parameters, and BIC and AIC differ only in their parameter number penalty, their outcomes were largely identical hence I used only the more stringent BIC. The model comparison metrics used were mean and median R^2 values, mean and median BIC values, and three others I designed using BIC values: *Win Count*, *Score* and *Total*. The *Win Count* finds which model generated the fit with the lowest BIC value for each timeseries and tallies them. Because that fails to consider how close the other BIC values were to this lowest value, the *Score*, derived from the aforementioned relative BIC value interpretation, assigns: five points to the winner and models within a BIC value of two; three to those within six; and one to those within ten, before summing them for each model. Because that too fails to consider some models may have performed significantly worse than the winner, i.e. a BIC value of ten to multiple hundreds higher, the *Total* for each timeseries calculates the total of the difference between each model and all models it outperformed, prior to summing them for each model.

3 Results

3.1 Model Fitting

After excluding the small datasets, 290 timeseries remained. For each of these, fits were successfully generated with all eight models, save for the Buchanan model failing to fit three. However, a successful fit does not imply a close fit, as demonstrated by **Figure 2**.

3.2 Model Comparison

The results of analyses comparing relative model performance, outlined in Section 2.4, are shown in **Table 1**. The quartic model achieved the best score in all categories, following a clear correlation of higher-order LMs performing better across the board. The logistic model performed similarly to the quadratic, while the other NLMs placed between the cubic and quartic LMs.

Similar analyses comparing only the four NLMs, shown in **Table 2**, prevent potentially higher-performing but biologically unsubstantial LMs from masking the results of the more meaningful NLMs. Here the Gompertz model is the clear winner, top in six out of seven metrics, narrowly losing to Baranyi on mean BIC value. Baranyi performed similarly, drawing with Gompertz on win count. The Buchanan model performed worse on all counts than Baranyi, as did the logistic model, trailing significantly, compared to the Buchanan model.

Table 1: Results of analysis comparing the fits for all 290 timeseries produced by each model

| Model | Mean R^2 | Median R^2 | Mean BIC | Median BIC | Win Count | Score | Total | Tally |
|-----------|------------|--------------|----------|------------|-----------|-------|--------|-------|
| Linear | 0.7759 | 0.859 | 21.73 | 25.76 | 3 | 76 | 2,179 | 0 |
| Quadratic | 0.9093 | 0.954 | 6.31 | 18.53 | 27 | 235 | 10,570 | 0 |
| Cubic | 0.9495 | 0.986 | -2.14 | 10.56 | 51 | 461 | 17,940 | 0 |
| Quartic | 0.9683 | 0.992 | -8.2 | 6.75 | 115 | 806 | 26,530 | 7 |
| Logistic | 0.7746 | 0.939 | 11.21 | 22.12 | 34 | 237 | 10,750 | 0 |
| Gompertz | 0.9542 | 0.99 | -3.94 | 8.78 | 97 | 639 | 20,980 | 0 |
| Baranyi | 0.953 | 0.9885 | -3.15 | 8.25 | 71 | 572 | 18,920 | 0 |
| Buchanan | 0.8221 | 0.987 | -2.28 | 11.1 | 67 | 497 | 18,660 | 0 |

Table 2: Results of analysis comparing the fits for all 290 timeseries produced by each NLM

| Model | Mean R^2 | Median R^2 | Mean BIC | Median BIC | Win Count | Score | Total | Tally |
|----------|------------|--------------|----------|------------|-----------|-------|-------|-------|
| Logistic | 0.7746 | 0.939 | 11.21 | 22.12 | 63 | 402 | 2,172 | 0 |
| Gompertz | 0.9542 | 0.99 | -3.94 | 8.78 | 137 | 898 | 7,055 | 6 |
| Baranyi | 0.953 | 0.9885 | -3.15 | 8.25 | 137 | 866 | 6,595 | 2 |
| Buchanan | 0.8221 | 0.987 | -2.28 | 11.1 | 109 | 744 | 5,925 | 0 |

3.3 Model-Covariable Correlation

The first of two figures reported here for the visualisation of patterns relating to the covariables of each experiment is **Figure 3**. The method behind it seeks to reveal if certain experimental conditions systematically affect the morphology of bacterial growth curves generated by one of the three highest-performing NLMs and, by extension, bacterial growth behaviours.

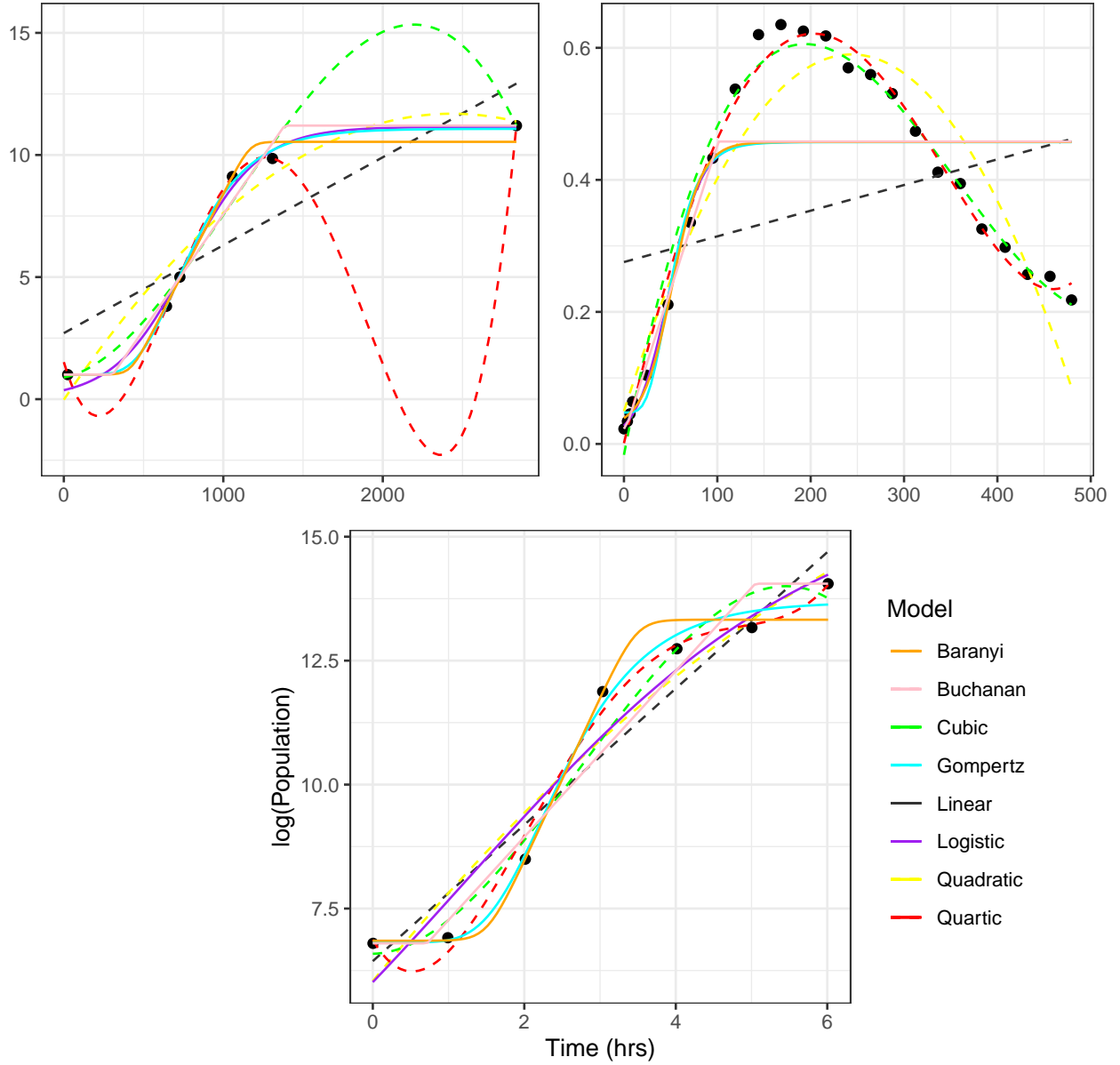


Figure 2: Exemplary timeseries plotted with the regression line that each model fit to them. Plots showing \log_2 of the population measurements against time in hours, which demonstrate to what extent each model can tolerate peculiar datasets. Linear models fit datasets comprising very few datapoints with near-perfect accuracy, yet without describing the true relationship at all (top left). Timeseries with a death phase and/or no lag phase are modelled poorly by the NLMs which plateau prematurely, although the logistic model benefits from a lack of lag phase (top right). For timeseries with lag phases but which had not plateaued when measurement ceased due to a drawn-out transition between exponential and stationary phases, logistic fails to capture the lag phase and Buchanan simply plateaus at the final datapoint, while Baranyi and Gompertz plateau harshly at the start of the transition, the latter to a lesser extent (bottom).

In contrast, **Figure 4** is the visualisation of an attempt to correlate individual model performance with different covariable categories. The idea is to suggest if one NLM may be more appropriate than another when modelling experimental data collected under certain conditions.

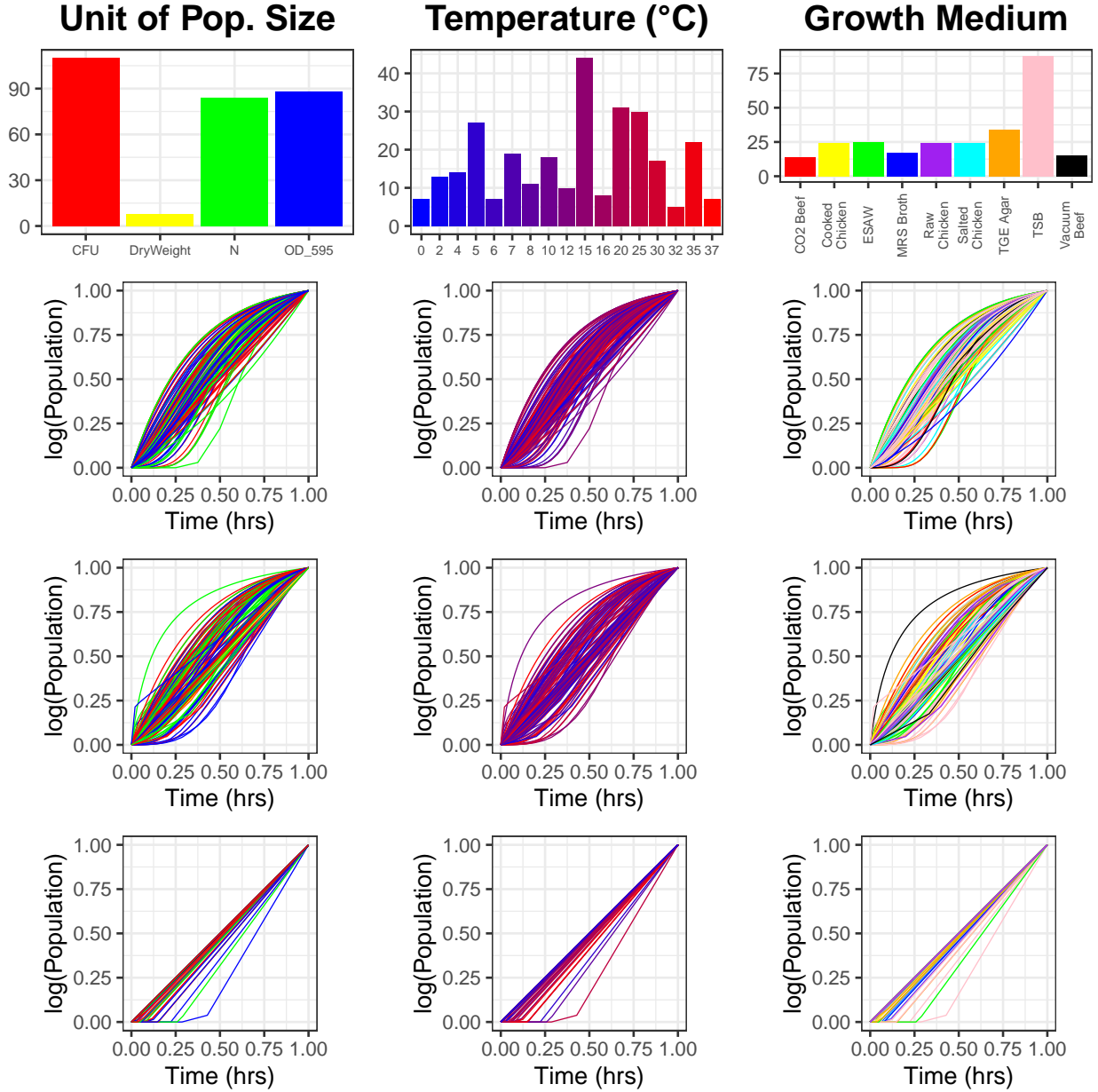


Figure 3: Standardised, superimposed growth curves for the Gompertz, Baranyi and Buchanan models, coloured by covariable category. Each fit of the three models was standardised in several steps. Firstly, lag phases and plateaus were removed, estimated as t_{lag} and 95% of N_{max} respectively. The resulting curves were transformed to start at the origin, all population values were divided by the highest remaining population value, and all time values by the highest remaining time value. Three copies of these normalised superimposed regression lines of Gompertz, Baranyi and Buchanan (rows 2-4 respectively) are coloured by categories of the three covariables. The barcharts displayed in row 1 show the abundancies of each category out of the 256 datasets, the bar colours corresponding to the regression lines beneath them.

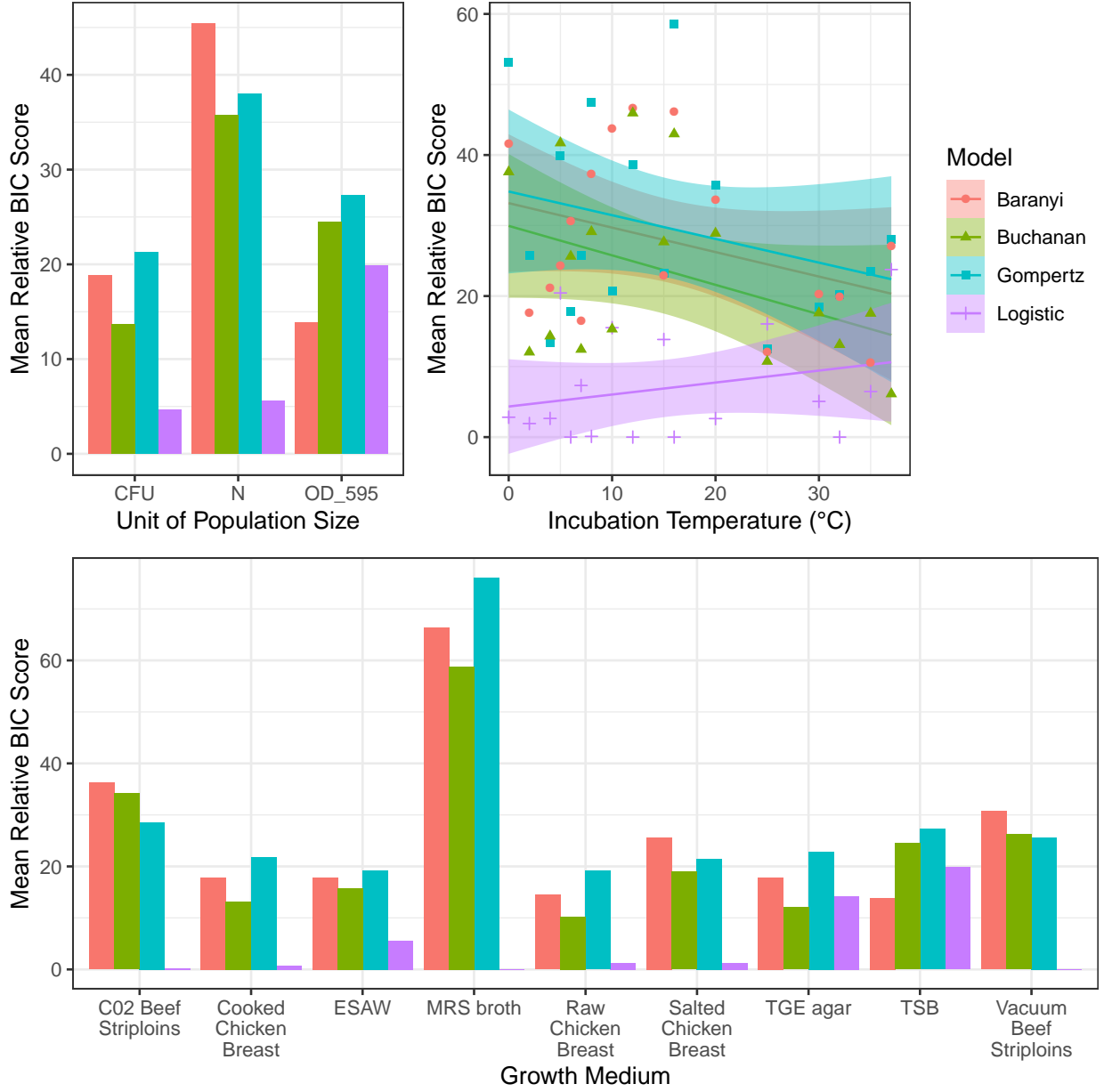


Figure 4: Plots expressing the relationships between the mean relative BIC scores of the four NLMs and the covariable categories. Mean relative BIC score here is the *Total* score, as defined in Section 2.4, awarded to each model for each timeseries, averaged across each covariable category. Barplots are used for the categoric covariables, while a scatterplot with linear regression lines and standard error ribbons is used to plot the continuous incubation temperature data. The colours of the legend apply to all three plots.

4 Discussion

I set out to devise a methodology to fit multiple models to hundreds of datasets, compare the performance of these models, and visualise biases certain models may have for certain experimental conditions.

The model fitting algorithm was highly effective, but illuminated the shortcomings of all eight models (**Figure 2**). The four NLMs are restrained by their parameters from fitting growth curves deviating from the standard lag, exponential and stationary phases. The logistic model, lacking the t_{lag} parameter, was additionally intolerant of lag phases but performed well absent one; the converse is true for the others. Gompertz best modelled atypical exponential phases, but all failed to capture death phases.

Meanwhile LMs are more malleable to irregular shapes, but restrained by their fixed number of parameters and innate curvature. One way to mitigate this restraint is with splines: linear models which produce smooth curves to fit the data.(White, 2017) As the linear regressions used herein rely upon OLS, they mathematically mould to the data, and will fit it perfectly where the dataset size equals one plus the number of linear model parameters; effectively "connecting the dots". This is known as Runge's phenomenon.(White, 2017) If there are other datapoints conveniently lined up as is common with empirical data of established natural relationships like population growth, LMs will also closely fit larger datasets with closeness of fit increasing with the order of linear model as **Table 1** demonstrates. This accounts for this deceitfully superior quartic model result. Its aforementioned weaknesses, which allow the NLMs to outperform it in typical growth curves, are outweighed by its superior but phenomenological, biologically inane fits of atypical ones: a shortfall that cannot be resolved with splines.

While this is in stark contrast with the theory-based NLMs fit using Levenberg-Marquardt non-linear least squares, an iterative method, these too have systematic problems which need resolution. Primarily, none of them address the death phase. Others have attempted to incorporate it into models but absent scientific basis, rendering them as uninformative as the dismissed LMs.(Peleg, 1997) Furthermore, they are innately inflexible with incomplete datasets missing one or both peripheral population growth phases. Ideally an alteration could be made to their equations enabling them to act similarly to splines, which allow the linear modelling of a supplied dataset while allowing for there to be an implicit 'bigger picture' behind the data not explicitly evident. While we could reactively resort to case-by-case altering of the model, such as how a simplified Baranyi model can be used to model growth without a plateau, that method is incompatible with the high-throughput focus of this paper.(Baranyi et al., 1993) Finally, they are sufficient as generic models for typical growth curves, but to accurately predict bacterial growth in fundamentally different scenarios a series of far more complex models which can be scientifically related to bacterial behaviour, physiology and biochemistry would be needed.

Model comparison heralded the modified Gompertz model as the most consistently high-performing NLM. This is corroborated by the conclusions of similar papers, which in turn supports the designed model comparison methodology.(R. L. Buchanan et al., 1997; Zwietering et al., 1990)

While the workflow in its current form is highly specific to the dataset and the models, its infrastructure was designed in some instances for more general use. Hence this workflow could in theory be repurposed into a bioinformatics pipeline to explore other groups of models, provided they can be translated into R functions, they model the same relationship, and the dataset both follows that relationship and includes potential covariate data.

The efficacy of the visualisation of covariables was less conclusive. In **Figure 3**, the only visible pattern is the Buchanan model suggesting growth rate increases with temperature. As for **Figure 4**, one must consider sample size before drawing conclusions, however the Baranyi model seems the superior choice when measuring population size with the "N" method explained in Section 2.2.1, whereas one may fare better with Gompertz when modelling bacterial growth on TSB. These are visual patterns with no formal statistical backing, but the objective was always to point towards potential lines of investigation, and could continue to with alternative groups of models.

An extra step not undertaken here is to perform similar visualisation on the poor datasets filtered out after fitting, to check for correlations between the other variables and highly irregular growth curves. Similarly, the visualisation would be more meaningful with larger and more equal samples, the collections of experiments with, for example, the same growth medium.

Ultimately, improving such visualisation could inform not only further investigation but also model development. Covariable regression could allow growth rate to be defined additionally by temperature, medium and so forth; after all, real-world bacterial growth is anisothermic. This could potentially be achieved by systematically incorporating dummy variables into a model using a base dataset like the one used herein, and observing if the model is significantly improved by it when tested on further datasets.(Suits, 1957)

In summation the Gompertz model, with Baranyi close behind, appear the best of those tested, but the parameters with which they are hard-coded mean they inherently require both a lag and stationary phase. This is true even where the timeseries in question would have been a typical sigmoidal growth curve had additional datapoints either side been collected. Moreover, this result is consistent with previous findings, lending credence to the model comparison algorithm used herein. While this did not discount the quartic model, it is widely accepted such LMs are tools for local prediction and not for extrapolation or mechanistic modelling. The overall methodology proved effective, and with the recommended alterations a repurposed version of it could enable the comparison of other models and the formation of new hypotheses.

References

- Akaike, H. (1974). A New Look at the Statistical Model Identification. *IEEE Transactions on Automatic Control*, 19(6), 716–723.
- Al-qadiri, H. M., Al-alamy, N. I., Lin, M., Al-holy, M., Cavinato, A. G., & Rasco, B. a. (2008). Fourier Transform Infrared Spectroscopy and Multivariate Analysis. *Journal of Rapid Methods & Automation in Microbiology*, 16, 73–89.
- Baranyi, J., Roberts, T. A., & McClure, P. (1993). A non-autonomous differential equation to model bacterial growth. *Food Microbiology*, 10(1), 43–59.
- Bolker, B. M., Gardner, B., Maunder, M., Berg, C. W., Brooks, M., Comita, L., Crone, E., Cubaynes, S., Davies, T., de Valpine, P., Ford, J., Gimenez, O., Kéry, M., Kim, E. J., Lennert-Cody, C., Magnusson, A., Martell, S., Nash, J., Nielsen, A., . . . Zipkin, E. (2013). Strategies for fitting nonlinear ecological models in R, AD Model Builder, and BUGS. *Methods in Ecology and Evolution*, 4(6), 501–512.
- Buchanan, R. E. (1918). Life Phases in a Bacterial Culture. *The Journal of Infectious Diseases*, 23(2), 109–125.
- Buchanan, R. L., Whiting, R. C., & Damert, W. C. (1997). A Comparison of the Gompertz, Baranyi, and three-phase linear models for fitting bacterial growth curves, 313–326.
- Grijspeerdt, K., & Vanrolleghem, P. (1999). Estimating the parameters of the Baranyi model for bacterial growth. *Food Microbiology*, 16(6), 593–605.
- Johnson, J. B., & Omland, K. S. (2004). Model selection in ecology and evolution. *Trends in Ecology and Evolution*, 19(2), 101–108.
- Kingsland, S. (1982). The Refractory Model : The Logistic Curve and the History of Population Ecology. *The Quarterly Review of Biology*, 57(1), 29–52.
- Mokhtari, M. S., Borzi, N. K., Foz, M. A., & Behzadi, M. R. (2019). Evaluation of non-linear models for genetic parameters estimation of growth curve traits in Kermani sheep. *Tropical Animal Health and Production*, 51(8), 2203–2212.
- Odenbaugh, J. (2006). The strategy of model building in population biology. *Biology and Philosophy*, 21(5), 607–621.
- Ozgul, A., Childs, D. Z., Oli, M. K., Armitage, K. B., Blumstein, D. T., Olson, L. E., Tuljapurkar, S., & Coulson, T. (2010). Coupled dynamics of body mass and population growth in response to environmental change. *Nature*, 466(7305), 482–485.
- Peleg, M. (1997). Modeling Microbial Populations with the Original and Modified Versions of the Continuous and Discrete Logistic Equations. *Critical Reviews in Food Science and Nutrition*, 37(5), 471–490.
- Peleg, M., & Corradini, M. G. (2011). Microbial growth curves: What the models tell us and what they cannot. *Critical Reviews in Food Science and Nutrition*, 51(10), 917–945.
- Posada, D., & Buckley, T. R. (2004). Model selection and model averaging in phylogenetics: Advantages of akaike information criterion and bayesian approaches over likelihood ratio tests. *Systematic Biology*, 53(5), 793–808.

317 Suits, D. B. (1957). Use of Dummy Variables in Regression Equations Author (s): Daniel B . Suits
318 Source : Journal of the American Statistical Association , Vol . 52 , No . 280 (Dec ., 1957
319), pp . Published by : Taylor & Francis , Ltd . on behalf of the American Statistica. *Journal*
320 *of the American Statistical Association*, 52(280), 548–551.

321 Tjørve, K. M., & Tjørve, E. (2017). The use of Gompertz models in growth analyses, and new
322 Gompertz-model approach: An addition to the Unified-Richards family. *PLoS ONE*, 12(6),
323 1–17.

324 Transtrum, M. K., & Qiu, P. (2016). Bridging Mechanistic and Phenomenological Models of Com-
325 plex Biological Systems. *PLoS Computational Biology*, 12(5)arXiv 1509.06278, 1–34.

326 Vrieze, S. I. (2012). Model selection and psychological theory: A discussion of the differences be-
327 tween the Akaike information criterion (AIC) and the Bayesian information criterion (BIC).
328 *Psychological Methods*, 17(2), 228–243.

329 White, D. R. (2017). Propagation of Uncertainty and Comparison of Interpolation Schemes. *Inter-*
330 *national Journal of Thermophysics*, 38(3), 1–14.

331 Zwietering, M., Jongenburger, I., Rombouts, F. M., & Van 'Triet, K. (1990). Modeling of the
332 Bacterial Growth Curve. *Applied and Environmental Microbiology*, 56(6), 1875–1881.