



Using generalized additive models to analyze biomedical non-linear longitudinal data

Beyond repeated measures ANOVA and Linear Mixed Models

Ariel I. Mundo , Timothy J. Muldoon*

Department of Biomedical Engineering, University of Arkansas, Fayetteville, AR, USA

tmuldoon@uark.edu

John R. Tipton 

Department of Mathematical Sciences, University of Arkansas, Fayetteville, AR, USA

1 Summary

In biomedical research, the outcome of longitudinal studies has been traditionally analyzed using the *repeated measures analysis of variance* (rm-ANOVA) or more recently, *linear mixed models* (LMEMs). Although LMEMs are less restrictive than rm-ANOVA in terms of correlation and missing observations, both methodologies share an assumption of linearity in the measured response, which results in biased estimates and unreliable inference when they are used to analyze data where the trends are non-linear, which is a common occurrence in biomedical research.

In contrast, *generalized additive models* (GAMs) relax the linearity assumption, and allow the data to determine the fit of the model while permitting missing observations and different correlation structures. Therefore, GAMs present an excellent choice to analyze non-linear longitudinal data in the context of biomedical research. This paper summarizes the limitations of rm-ANOVA and LMEMs and uses simulated data to visually show how both methods produce biased estimates when used on non-linear data. We also present the basic theory of GAMs, and using trends of oxygen saturation in tumors reported in the biomedical literature, we simulate example longitudinal data (2 treatment groups, 10 subjects per group, 5 repeated measures for each group) to demonstrate how these models can be computationally implemented. We show that GAMs are able to produce estimates that are consistent with the trends of biomedical non-linear data even in the case when missing observations exist (with 40% of the simulated observations missing), allowing reliable inference from the data. To make this work reproducible, the code and data used in this paper are available at: <https://github.com/aimundo/GAMs-biomedical-research>.

Keywords

cancer biology; tumor response; generalized additive models; simulation; R

2 Background

Longitudinal studies are designed to repeatedly measure a variable of interest in a group (or groups) of subjects, with the intention of observing the evolution of effect across time rather than analyzing a single time point (e.g., a cross-sectional study). Biomedical research frequently uses longitudinal studies to analyze the evolution of a “treatment” effect across multiple time points; and in such studies the subjects of analysis range from animals (mice, rats, rabbits), to human patients, cells, or blood samples, among many others. Tumor response (Roblyer *and others* 2011; Demidov *and others* 2018; Pavlov *and others* 2018; Tank *and others* 2020), antibody expression (Ritter *and others* 2001; Roth *and others* 2017), and cell metabolism (Skala *and others* 2010; Jones *and others* 2018) are examples of the different situations where researchers have used longitudinal designs to study some physiological response. Because the frequency of the measurements in a longitudinal study is dependent on the biological phenomena of interest and the experimental design of the study, the frequency of such measurements can range from minute intervals to study a short-term response such as anesthesia effects in animals (Greening *and others* 2018), to weekly measurements to analyze a mid-term response like the evolution of dermatitis symptoms in breast cancer patients (Sio *and others* 2016), to monthly measurements to study a long-term response such as mouth opening following radiotherapy (RT) in neck cancer patients (Kamstra *and others* 2015).

Traditionally, a “frequentist” or “classical” statistical paradigm is used in biomedical research to derive inferences from a longitudinal study. The frequentist paradigm regards probability as the limit of the expected outcome when an experiment is repeated a large number of times (Wagenmakers *and others* 2008), and such view is applied to the analysis of longitudinal data by assuming a null hypothesis under a statistical model that is often an *analysis of variance over repeated measures* (repeated measures ANOVA or rm-ANOVA). The rm-ANOVA model makes three key assumptions regarding longitudinal data: 1) linearity of the response across time, 2) constant correlation across same-subject measurements, and 3) observations from each subject are obtained at all time points through the study (a condition also known as *complete observations*) (Gueorguieva and Krystal 2004; Schober and Vetter 2018).

The expected linear behavior of the response through time is a key requisite in rm-ANOVA (Pinheiro and Bates 2006). This “linearity assumption” in rm-ANOVA implies that the model is misspecified when the data does not follow a linear trend, which results in unreliable inference. In biomedical research, non-linear

trends are the norm rather than the exception in longitudinal studies. A particular example of this non-linear behavior in longitudinal data arises in measurements of tumor response to chemo and/or radiotherapy in preclinical and clinical settings (Vishwanath *and others* 2009; Skala *and others* 2010; Roblyer *and others* 2011). These studies have shown that the collected signal does not follow a linear trend over time, and presents extreme variability at different time points, making the fit of rm-ANOVA model inconsistent with the observed variation. Therefore, when rm-ANOVA is used to draw inference of such data the estimates are inevitably biased, because the model is only able to accommodate linear trends that fail to adequately represent the biological phenomenon of interest.

A *post hoc* analysis is often used in conjunction with rm-ANOVA to perform repeated comparisons to estimate a *p-value*, which in turn is used as a measure of significance. Although it is possible that a *post hoc* analysis of rm-ANOVA is able to find “significant” *p-values* ($p < 0.05$) from non-linear data, the validity of such metric is dependent on how adequate the model fits the data. In other words, *p-values* are valid only if the model and the data have good agreement; if that is not the case, a “Type III” error (known as “model misspecification”) occurs (Dennis *and others* 2019). For example, model misspecification will occur when a model that is only able to explain linear responses (such as rm-ANOVA) is fitted to data that follows a quadratic trend, thereby causing the resulting *p-values* and parameter estimates to be invalid (Wang *and others* 2019).

Additionally, the *p-value* itself is highly variable, and multiple comparisons can inflate the false positivity rate (Type I error or α) (Liu *and others* 2010; Halsey *and others* 2015), consequently biasing the conclusions of the study. Corrections exist to address the Type I error issue of multiple comparisons (such as Bonferroni (Abdi 2010)), but they in turn reduce statistical power ($1 - \beta$) (Nakagawa 2004), and lead to increased Type II error (failing to reject the null hypothesis when the null hypothesis is false) (Gelman *and others* 2012; Albers 2019). Therefore, the tradeoff of *post hoc* comparisons in rm-ANOVA between Type I, II and III errors might be difficult to resolve in a biomedical longitudinal study where a delicate balance exists between statistical power and sample size.

On the other hand, the assumption of constant correlation in rm-ANOVA (often known as the *compound symmetry assumption*) is typically unreasonable because correlation between the measured responses often diminishes as the time interval between the observation increases (Ugrinowitsch *and others* 2004). Corrections can be made in rm-ANOVA in the absence of compound symmetry (Greenhouse and Geisser 1959; Huynh and Feldt 1976), but the effectiveness of the correction is limited by the size of the sample, the number of measurements (Haverkamp and Beauducel 2017), and group sizes (Keselman *and others* 2001). In the case of biomedical research, where living subjects are frequently used, sample sizes are often not “large” due to ethical and budgetary reasons (Charan and Kantharia 2013) which might cause the corrections for lack of

compound symmetry to be ineffective.

Due to a variety of causes, the number of observations during a study can vary between all subjects. For example, in a clinical trial patients may voluntarily withdraw, whereas attrition due to injury or weight loss in preclinical animal studies is possible. It is even plausible that unexpected complications with equipment or supplies arise that prevent the researcher from collecting measurements at certain time points. In each of these missing data scenarios, the *complete observations* assumption of classical rm-ANOVA is violated. When incomplete observations occur, a rm-ANOVA model is fit by excluding all subjects with missing observations from the analysis (Gueorguieva and Krystal 2004). This elimination of partially missing data from the analysis can result in increased costs if the desired statistical power is not met with the remaining observations, because it would be necessary to enroll more subjects. At the same time, if the excluded observations contain insightful information that is not used, their elimination from the analysis may limit the demonstration of significant differences between groups.

During the last decade, the biomedical community has started to recognize the limitations of rm-ANOVA in the analysis of longitudinal data. The recognition on the shortcomings of rm-ANOVA is exemplified by the use of linear mixed effects models (LMEMs) by certain groups to analyze longitudinal tumor response data (Vishwanath *and others* 2009; Skala *and others* 2010). Briefly, LMEMs incorporate *fixed effects*, which correspond to the levels of experimental factors in the study (e.g., the different drug regimens in a clinical trial), and *random effects*, which account for random variation within the population (e.g., the individual-level differences not due to treatment such as weight or age). When compared to the traditional rm-ANOVA, LMEMs are more flexible as they can accommodate missing observations for multiple subjects and allow different modeling strategies for the variability within each measure in every subject (Pinheiro and Bates 2006). However, LMEMs impose restrictions in the distribution of the errors of the random effects, which need to be normally distributed and independent (Gueorguieva and Krystal 2004; Barr *and others* 2013). And even more importantly, LMEMs also assume a linear relationship between the response and time (Pinheiro and Bates 2006), making them unsuitable to analyze non-linear data.

As the rm-ANOVA and the more flexible LMEM approaches make overly restrictive assumptions regarding the linearity of the response, there is a need for biomedical researchers to explore the use of additional statistical tools that allow the data (and not an assumption in trend) to determine the trend of the fitted model, to enable appropriate inference. In this regard, generalized additive models (GAMs) present an alternative approach to analyze longitudinal data. Although not frequently used by the biomedical community, these semi-parametric models are customarily used in other fields to analyze longitudinal data. Examples of the use of GAMs include the analysis of temporal variations in geochemical and palaeoecological data

(Rose *and others* 2012; Simpson 2018; Pedersen *and others* 2019), health-environment interactions (Yang *and others* 2012) and the dynamics of government in political science (Beck and Jackman 1998). There are several advantages of GAMs over LMEMs and rm-ANOVA models: 1) GAMs can fit a more flexible class of smooth responses that enable the data to dictate the trend in the fit of the model, 2) they can model non-constant correlation between repeated measurements (Wood 2017) and 3) can easily accommodate missing observations. Therefore, GAMs can provide a more flexible statistical approach to analyze non-linear biomedical longitudinal data than LMEMs and rm-ANOVA.

The current advances in programming languages designed for statistical analysis (specifically R), have eased the computational implementation of traditional models such as rm-ANOVA and more complex approaches such as LMEMs and GAMs. In particular, R (R Core Team 2020) has an extensive collection of documentation and functions to fit GAMs in the package *mgcv* (Wood *and others* 2016; Wood 2017) that not only speed up the initial stages of the analysis but also enable the use of advanced modeling structures (e.g. hierarchical models, confidence interval comparisons) without requiring advanced programming skills from the user. At the same time, R has many tools that simplify data simulation, an emerging strategy used to test statistical models (Haverkamp and Beauducel 2017). Data simulation methods allow the researcher to create and explore different alternatives for analysis without collecting information in the field, reducing the time window between experiment design and its implementation, and simulation can be also used for power calculations and study design questions.

This work provides biomedical researchers with a clear understanding of the theory and the practice of using GAMs to analyze longitudinal data using by focusing on four areas. First, the limitations of LMEMs and rm-ANOVA regarding linearity of response, constant correlation structures and missing observations are explained in detail. Second, the key theoretical elements of GAMs are presented using clear and simple mathematical notation while explaining the context and interpretation of the equations. Third, we illustrate the type of non-linear longitudinal data that often occurs in biomedical research using simulated data that reproduces patterns in previously reported studies (Vishwanath *and others* 2009). The simulated data experiments highlight the differences in inference between rm-ANOVA, LMEMs and GAMs on data similar to what is commonly observed in biomedical studies. Finally, reproducibility is emphasized by providing the code to generate the simulated data and the implementation of different models in R, in conjunction with a step-by-step guide demonstrating how to fit models of increasing complexity.

In summary, this work will allow biomedical researchers to identify when the use of GAMs instead of rm-ANOVA or LMEMs is appropriate to analyze longitudinal data, and provide guidance on the implementation of these models to improve the standards for reproducibility in biomedical research.

3 Challenges presented by longitudinal studies

3.1 The repeated measures ANOVA and Linear Mixed Model

The *repeated measures analysis of variance* (rm-ANOVA) and the *linear mixed model* (LMEM) are the most commonly used statistical analysis for longitudinal data in biomedical research. These statistical methodologies require certain assumptions for the model to be valid. From a practical view, the assumptions can be divided in three areas: 1) linear relationship between covariates and response, 2) a constant correlation between measurements, and, 3) complete observations for all subjects. Each one of these assumptions is discussed below.

3.2 Linear relationship

3.2.1 The repeated measures ANOVA case

In a longitudinal biomedical study, two or more groups of subjects (e.g., human subject, mice, samples) are subject to different treatments (e.g., a “treatment” group receives a novel drug or intervention vs. a “control” group that receives a placebo), and measurements from each subject within each group are collected at specific time points. The collected response is modeled with *fixed* components. The *fixed* component can be understood as a constant value in the response which the researcher is interested in measuring, i.e., the average effect of the novel drug/intervention in the “treatment” group.

Mathematically speaking, a rm-ANOVA model with an interaction can be written as:

$$y_{ijt} = \beta_0 + \beta_1 \times time_t + \beta_2 \times treatment_j + \beta_3 \times time_t \times treatment_j + \varepsilon_{ijt} \quad (1)$$

In this model y_{ijt} is the response for subject i , in treatment group j at time t , which can be decomposed in a mean value β_0 , *fixed effects* of time ($time_t$), treatment ($treatment_j$) and their interaction $time_t * treatment_j$ which have linear slopes given by β_1, β_2 and β_3 , respectively. Independent errors ε_{ijt} represent random variation not explained by the *fixed* effects, and are assumed to be $\sim N(0, \sigma^2)$ (independently and identically normally distributed with mean zero and variance σ^2). In a biomedical research context, suppose two treatments groups are used in a study (e.g., “placebo” vs. “novel drug” or “saline” vs. “chemotherapy”). Then, the group terms in Equation (1) can be written as below with $treatment_j = 0$ representing the first treatment group (Group A) and $treatment_j = 1$ representing the second treatment group (Group B). With this notation, the linear model then can be expressed as

$$y_{ijt} = \begin{cases} \beta_0 + \beta_1 \times time_t + \varepsilon_{ijt} & \text{if Group A} \\ \beta_0 + \beta_2 + \beta_1 \times time_t + \beta_3 \times time_t + \varepsilon_{ijt} & \text{if Group B} \end{cases} \quad (2)$$

180 To further simplify the expression, substitute $\widetilde{\beta}_0 = \beta_0 + \beta_2$ and $\widetilde{\beta}_1 = \beta_1 + \beta_3$ in the equation for Group B.
 181 This substitution allows for a different intercept and slope for Groups A and B. The model is then written as

$$y_{ijt} = \begin{cases} \beta_0 + \beta_1 \times time_t + \varepsilon_{ijt} & \text{if Group A} \\ \widetilde{\beta}_0 + \widetilde{\beta}_1 \times time_t + \varepsilon_{ijt} & \text{if Group B} \end{cases} \quad (3)$$

182 Presenting the model in this manner makes clear that when treating different groups, an rm-ANOVA model
 183 is able to accommodate non-parallel lines in each case (different intercepts and slopes per group). In other
 184 words, the rm-ANOVA model “expects” a linear relationship between the covariates and the response, this
 185 means that either presented as Equation (1), Equation (2) or Equation (3), an rm-ANOVA model is only able
 186 to accommodate linear patterns in the data. If the data show non-linear behavior, the rm-ANOVA model
 187 will approximate this behavior with non-parallel lines.

188 3.2.2 The Linear Mixed Model Case (LMEM)

189 A LMEM is a class of statistical models that incorporates *fixed effects* to model the relationship between the
 190 covariates and the response, and *random effects* to model subject variability that is not the primary focus of
 191 the study but that might be important to distinguish (Pinheiro and Bates 2006; West *and others* 2014). A
 192 LMEM with interaction between time and treatment for a longitudinal study can be written as:

$$y_{ijt} = \beta_0 + \beta_1 \times time_t + \beta_2 \times treatment_j + \beta_3 \times time_t \times treatment_j + \mu_{ij} + \varepsilon_{ijt} \quad (4)$$

193 When Equation (1) and Equation (4) are compared, it is easily noticeable that LMEM and rm-ANOVA have
 194 the same construction regarding the *fixed effects* of time and treatment, but that the LMEM incorporates an
 195 additional source of variation (the term μ_{ij}). This term μ_{ij} is the one that corresponds to the *random effect*,
 196 accounting for variability in each subject (subject_i) within each group (group_j). The *random* component can
 197 also be understood as used to model some “noise” in the response, but that is intended to be analyzed and
 198 disentangled from the “global noise” term ε_{ijt} from Equation (1).

199 For example, if the blood concentration of the drug is measured in certain subjects in the early hours of

the morning while other subjects are measured in the afternoon, it is possible that the difference in the collection time introduces some “noise” in the data. As the name suggests, this “random” variability needs to be modeled as a variable rather than as a constant value. The *random effect* μ_{ij} in Equation (4) is assumed to be $\mu_{ij} \sim N(0, \sigma_\mu^2)$. In essence, the *random effect* in a LMEM enables to fit models with different slopes at the subject-level (Pinheiro and Bates 2006). However, the expected linear relationship of the covariates and the response in Equation (1) and in Equation (4) is essentially the same, representing a major limitation of LMEMs to fit a non-linear response.

3.3 Covariance in rm-ANOVA and LMEMs

In a longitudinal study there is an expected *covariance* between repeated measurements on the same subject, and because repeated measures occur in the subjects within each group, there is a *covariance* between measurements at each time point within each group. The *covariance matrix* (also known as the variance-covariance matrix) is a matrix that captures the variation between and within subjects in a longitudinal study (Wolfinger 1996) (For an in-depth analysis of the covariance matrix see (Weiss 2005; West *and others* 2014)).

In the case of an rm-ANOVA analysis, it is typically assumed that the covariance matrix has a specific construction known as *compound symmetry* (also known as “sphericity” or “circularity”). Under this assumption, the between-subject variance and within-subject correlation are constant across time (Geisser and Greenhouse 1958; Huynh and Feldt 1976; Weiss 2005). However, it has been shown that this condition is frequently not justified because the correlation between measurements tends to change over time (Maxwell *and others* 2017); and it is higher between consecutive measurements (Gueorguieva and Krystal 2004; Ugrinowitsch *and others* 2004). Although corrections can be made (such as Huynh-Feldt or Greenhouse-Geisser) (Greenhouse and Geisser 1959; Huynh and Feldt 1976) the effectiveness of each correction is limited because it depends on the size of the sample, the number of repeated measurements (Haverkamp and Beauducel 2017), and they are not robust if the group sizes are unbalanced (Keselman *and others* 2001). Because biomedical longitudinal studies are often limited in sample size and can have an imbalanced design, the corrections required to use an rm-ANOVA model may not be able to provide a reasonable adjustment that makes the model valid.

In the case of LMEMs, one key advantage over rm-ANOVA is that they allow different structures for the variance-covariance matrix including exponential, autoregressive of order 1, rational quadratic and others (Pinheiro and Bates 2006). Nevertheless, the analysis required to determine an appropriate variance-covariance structure for the data can be a challenging process by itself. Overall, the spherical assumption for rm-ANOVA may not capture the natural variations of the correlation in the data, and can bias the inferences from the

analysis.

3.4 Missing observations

Missing observations are an issue that arises frequently in longitudinal studies. In biomedical research, this situation can be caused by reasons beyond the control of the investigator (Molenberghs 2004). Dropout from patients and attrition or injury in animals are among the reasons for missing observations. Statistically, missing information can be classified as *missing at random* (MAR), *missing completely at random* (MCAR), and *missing not at random* (MNAR) (Weiss 2005). In a MAR scenario, the pattern of the missing information is related to some variable in the data, but it is not related to the variable of interest (Scheffer 2002). If the data are MCAR, this means that the missingness is completely unrelated to the collected information (Potthoff and others 2006), and in the case of MNAR the missing values are dependent on their value.

An rm-ANOVA model assumes complete observations for all subjects, and therefore subjects with one or more missing observations are excluded from the analysis. This is inconvenient because the remaining subjects might not accurately represent the population, and statistical power is affected by this reduction in sample size (Ma and others 2012). In the case of LMEMs, inferences from the model are valid when missing observations in the data exist that are MAR or MCAR (West and others 2014). For example, if attrition occurs in all mice that had lower weights at the beginning of a chemotherapy response study, the missing data can be considered MAR because the missingness is unrelated to other variables of interest.

3.5 What do an rm-ANOVA and LMEM fit look like? A visual representation using simulated data

To visually demonstrate the limitations of rm-ANOVA and LMEMs for non-linear longitudinal data, this section presents a simulation experiment of a normally distributed response of two groups of 10 subjects each. An rm-ANOVA model (Equation (1)), and a LMEM (Equation (4)) are fitted to each group, using R (R Core Team 2020) and the package *nlme* (Pinheiro and others 2020).

Briefly, two cases for the mean responses for each group are considered: in the first case, the mean response in each group is a linear function over time with different intercepts and slopes; a negative slope is used for Group 1 and a positive slope is used for Group 2 (Figure 1A). In the second case, a second-degree polynomial (quadratic) function is used for the mean response per group: the quadratic function is concave down for Group 1 and it is concave up for Group 2 (Figure 1C). In both the linear and quadratic simulated data, the groups start with the same mean value at the first time point. This is intentional in order to simulate the

expected temporal evolution of some physiological quantity, which is typical in biomedical experiments where a strong non-linear trend is present.

Specifically, the rationale for the chosen linear and quadratic functions is the expectation that a measured response in two treatment groups is similar in the initial phase of the study, but as therapy progresses a divergence in the trend of the response indicates a treatment effect. In other words, Group 1 can be thought as a “Control” group and Group 2 as a “Treatment” group. From the mean response per group (linear or quadratic), the variability or “error” of individual responses within each group is simulated using a covariance matrix with compound symmetry (constant variance across time). Thus, the response per subject in both the linear and quadratic simulation corresponds to the mean response per group plus the error (Figure 1 B,D).

A more comprehensive exploration of the fit of rm-ANOVA and LMEMs for linear and non-linear longitudinal data appears in the Appendix (Figure A.1 and Figure A.2), where simulation with compound symmetry and independent errors (errors generated from a normal distribution that are not constant over time) and the plot of simulated errors, and fitted parameters is presented. We are aware that the simulated data used in this section present an extreme case that might not occur frequently in biomedical research, but they are used as a representation of the consequences of modeling non-linear data with a linear model such as rm-ANOVA or LMEMs. Of notice, in Section 6 we use simulated data that does follow reported trends in the biomedical literature to implement GAMs.

The simulation shows that the fits produced by the LMEM and the rm-ANOVA model are good for linear data, as the predictions for the mean response are reasonably close to the “truth” of the simulated data (Figure 1A). When the linearity and compound symmetry assumptions are met, the rm-ANOVA model approximates well the global trend by group (Figure 1B). Note that because the LMEM incorporates *random effects*, is able to provide estimates for each subject and a “global” estimate (Figure 1C).

However, consider the case when the data follows a non-linear trend, such as the simulated data in Figure 1D. Here, the mean response per group was simulated using a quadratic function, and errors and individual responses were produced as in Figure 1A. The mean response in the simulated data with quadratic behavior changes in each group through the timeline, and the mean value is the same as the initial value by the fifth time point for each group. Fitting an rm-ANOVA model (Equation (1)) or a LMEM (Equation (4)) to this data produces the fit that appears in Figure 1E, F.

Comparing the fitted responses of the LMEM and the rm-ANOVA models used in the simulated quadratic data (Figure 1E, F) indicates that the models are not capturing the changes within each group. Specifically, note that the fitted mean response of both models shows that the change (increase for Treatment 1 or decrease

for Treatment 2) in the response through time points 2 and 4 is not being captured. The LMEM is only able to account for between-subject variation by providing estimates for each subject (Figure 1F), but both models are unable to capture the fact that the initial values are the same in each group, and instead fit non-parallel lines that have initial values that are markedly different from the “true” initial values in each case (compare Figure 1D with Figure 1E, F). If such a change has important physiological implications, both rm-ANOVA and LMEMs omit it from the fitted mean response. Thus, even though the model correctly detects a divergence between treatment groups, the exact nature of this difference is not correctly identified, limiting valuable inferences from the data.

This section has used simulation to better convey the limitations of linearity and correlation in the response in non-linear data. The models fitted to the simulated data were an rm-ANOVA model and a LMEM, where the main issue is the expected linear trend in the response. In the following section, we present generalized additive models (GAMs) as a data-driven alternative method to analyze longitudinal non-linear data that overcomes the linearity assumption.

4 GAMs as a special case of Generalized Linear Models

4.1 GAMs and Basis Functions

Generalized linear models (GLMs) are a family of models (which include rm-ANOVA and LMEMs) that fit a linear response function to data that may not have normally distributed errors (Nelder and Wedderburn 1972). In contrast, GAMs are a family of regression-based methods for estimating smoothly varying trends and are a broader class of models that contain the GLM family as a special case (Hastie and Tibshirani 1987; Wood 2017; Simpson 2018). A GAM model can be written as:

$$y_{ijt} = \beta_0 + f(x_t | \beta_j) + \varepsilon_{ijt} \quad (5)$$

Where y_{ijt} is the response at time t of subject i in group j , β_0 is the expected value at time 0, the change of y_{ijt} over time is represented by the *smooth function* $f(x_t | \beta_j)$ with inputs as the covariates x_t and parameters β_j , and ε_{ijt} represents the residual error.

In contrast to the linear functions used to model the relationship between the covariates and the response in rm-ANOVA or LMEM, GAMs use more flexible *smooth functions*. This approach is advantageous as it does not restrict the model to a linear relationship, although a GAM can estimate a linear relationship if the data is consistent with a linear response. One possible set of functions for $f(x_t | \beta_j)$ that allow for non-linear

responses are polynomials, but a major limitation is that polynomials create a “global” fit as they assume that the same relationship exists everywhere, which can cause problems with inference (Beck and Jackman 1998). In particular, polynomial fits are known to show boundary effects because as t goes to $\pm\infty$, $f(x_t | \beta_j)$ goes to $\pm\infty$ which is almost always unrealistic and causes bias at the endpoints of the time period.

The smooth functional relationship between the covariates and the response in GAMs is specified using a semi-parametric relationship that can be fit within the GLM framework, by using *basis function* expansions of the covariates and by estimating random coefficients associated with these basis functions. A *basis* is a set of functions that spans the mathematical space where the smooths that approximate $f(x_t | \beta_j)$ exist (Simpson 2018). For the linear model in Equation (1), the basis coefficients are β_1 , β_2 and β_3 and the basis vectors are $time_t$, $treatment_j$ and $time_t \times treatment_j$. The basis function then, is the combination of basis coefficients and basis vectors that map the possible relationship between the covariates and the response (Hefley and others 2017), which in the case of Equation (1) is restricted to a linear family of functions. In the case of Equation (5), the basis functions are contained in the expression $f(x_t | \beta_j)$, which means that the model allows for non-linear relationships among the covariates.

Splines (cubic, thin plate, etc.) are commonly used *basis functions*; a cubic spline is a smooth curve constructed from cubic polynomials joined together in a manner that enforces smoothness, and thin plate regression splines are an optimized version that work well with noisy data (Wood 2017; Simpson 2018). Splines have a long history in solving semi-parametric statistical problems and are often a default choice to fit GAMs as they are a simple, flexible and powerful option to obtain smoothness (Wegman and Wright 1983). Therefore, this data-driven flexibility in GAMs overcomes the limitation that occurs in LMEMs and rm-ANOVA when the data is non linear.

To further clarify the concept of basis functions and smooth functions, consider the simulated response for Group 1 in Figure 1C. The simplest GAM model that can be used to estimate such response is that of a single smooth term for the time effect; i.e., a model that fits a smooth to the trend of the group through time. The timeline can be divided in equally spaced *knots*, each knot being a region where a different set of basis functions will be used. Because there are six timepoints for this group, five knots can be used. The model with five knots to construct the smooth term means that it will have four basis functions (plus one that corresponds to the intercept). The choice of basis functions is set using default values in the package *mgcv* depending on the number of knots. In Figure 2A, the four basis functions (and the intercept) are shown. Each of the basis functions is composed of six different points (because there are six points on the timeline). To control the “wiggleness” of the fit, each of the basis functions of Figure 2A is weighted by multiplying it by a coefficient according to the matrix of Figure 2B. The parameter estimates are penalized (shrunk towards 0)

where the penalty reduces the “wiggleness” of the smooth fit to prevent overfitting. A weak penalty estimate will result in wiggly functions whereas a strong penalty estimate provides evidence that a linear response is appropriate.

To get the weighted basis functions, each basis (from Figure Figure 2A) is multiplied by the corresponding coefficients in Figure 2B, thereby increasing or decreasing the original basis functions. Figure 2C shows the resulting weighted basis functions. Note that the magnitude of the weighting for the first basis function has resulted in a decrease of its overall value (because the coefficient for that basis function is less than 1). On the other hand, the third basis function has roughly doubled its value. Finally, the weighted basis functions are added at each timepoint to produce the smooth term. The resulting smooth term for the effect of *time* is shown in Figure 2D (orange line), along the simulated values per group, which appear as points.

5 A Bayesian interpretation of GAMs

Bayes’ theorem states that the probability of an event can be calculated using prior knowledge or belief (McElreath 2018). In the case of non-linear data, the belief that the *true* trend of the data is likely to be smooth rather than “wiggly” introduces the concept of a prior distribution for wiggleness (and therefore a Bayesian view) of GAMs (Wood 2017). GAMs are considered “empirical” Bayesian models because the smoothing parameters are estimated from the data (and not from a prior distribution as in the “Full Bayes” case) (Miller 2019). Moreover, the use of the restricted maximum likelihood (REML) to estimate the smoothing parameters gives an empirical estimate of the smooth model (Laird and Ware 1982; Pedersen *and others* 2019). Therefore, the confidence intervals calculated for the smooth terms using the package *mgcv* are considered empirical Bayesian posterior credible intervals (Pedersen *and others* 2019), which have good “frequentist” coverage (pointwise coverage or “single point” coverage), and *across the function* coverage (Wood 2017). This last part means that contrary to a pointwise coverage (where the coverage of the interval is correct for a single point) the estimated confidence intervals for the smooths will contain *on average* the true function of the data 95% of the time across the entire timeline (in the case of longitudinal data for which smooths are calculated), which allows to obtain better inference from the model. In-depth theory of the Bayesian interpretation of GAMs is beyond the scope of this paper, but can be found in (Wood 2017; Simpson 2018; Miller 2019) and (Marra and Wood 2012). With this brief introduction to the Bayesian interpretation of GAMs, we henceforth refer to the confidence intervals for the smooths in GAMs as “empirical Bayesian” through the rest of this paper.

6 The analysis of longitudinal biomedical data using GAMs

The previous sections provided the basic framework to understand the GAM framework and how these models are more advantageous to analyze non-linear longitudinal data when compared to rm-ANOVA or LMEMs. This section will use simulation to present the practical implementation of GAMs for longitudinal biomedical data using R and the package `mgcv`. The code for the simulated data and figures, and a brief guide for model selection and diagnostics appear in the Appendix.

6.1 Simulated data

The simulated data is based on the reported longitudinal changes in oxygen saturation (StO_2) in subcutaneous tumors that appear in Figure 3C in (Vishwanath *and others* 2009). In the paper, diffuse reflectance spectroscopy was used to quantify StO_2 changes in both groups at the same time points (days 0, 2, 5, 7 and 10). In the “Treatment” group (chemotherapy) an increase in StO_2 is observed through time, while a decrease is seen in the “Control” (saline) group. Following the reported trend, we simulated 10 normally distributed observations at each time point with a standard deviation (SD) of 10% (matching the SD in the original paper). The simulated and real data appear in Figure 3A and the inset, respectively.

6.2 An interaction GAM for longitudinal data

An interaction effect is typically the main interest in longitudinal biomedical data, as it takes into account treatment, time, and their combination. In a practical sense, when a GAM is implemented for longitudinal data, a smooth can be added to the model for the *time* effect to account for the repeated measures over time. Although specific methods of how GAMs model correlation structures is a topic beyond the scope of this paper, it suffices to say that GAMs are flexible and can handle correlation structures beyond compound symmetry. A detailed description on basis functions and correlations can be found in (Hefley *and others* 2017).

For the data in Figure 3, A the main effect of interest is how StO_2 changes over time for each treatment. To estimate this, the model incorporates independent smooths for *Group* and *Day*, respectively. The main thing to consider is that model syntax accounts for the fact that one of the variables is numeric (*Day*) and the other is a factor (*Group*). Because the smooths are centered at 0, the factor variable needs to be specified as a parametric term in order to identify any differences between the groups. Using R and the package `mgcv` the model syntax is:

```
m1 <- gam(StO2_sim ~ Group + s(Day, by=Group, k=5), method='REML', data =
```

```
dat_sim)
```

This syntax specifies that `m1` will store the model, and that the change in the simulated oxygen saturation (`StO2_sim`) is modeled using independent smooths over `Day` for each `Group` (the parenthesis preceded by `s`) using 5 knots. The smooth is constructed by default using thin plate regression splines, but other splines can be used if desired, including Gaussian process smooths (Simpson 2018). The parametric term `Group` is added to quantify overall mean differences in the effect of treatment between groups, and the `method` chosen to estimate the smoothing parameters is the restricted maximum likelihood (REML) (Wood 2017). When the smooths are plotted over the raw data, it is clear that the model has been able to capture the trend of the change of StO_2 for each group across time (Figure 3B). Model diagnostics can be obtained using the `gam.check` function, and the function `appraise` from the package *gratia* (Simpson 2020). A guide for model selection and diagnostics is in the Appendix, and an in-depth analysis can be found in (Wood 2017) and (Harezlak and others 2018).

One question that might arise at this point is “what is the fit that an rm-ANOVA model produces for the simulated data?” The rm-ANOVA model, which corresponds to Equation (1) is presented in Figure 3C. This is a typical case of model misspecification: The slopes of each group are different, which would lead to a *p-value* indicating significance for the treatment and time effects, but the model is not capturing the changes that occur at days 2 and between days 5 and 7, whereas the GAM model is able to reliably estimate the trend over all timepoints (Figure 3B) .

Because GAMs do not require equally-spaced or complete observations for all subjects, they are advantageous to analyze longitudinal data where missingness exists. The rationale behind this is that GAMs are able to pick the trend in the data even when some observations are missing. However, this usually causes the resulting smooths to have wider confidence intervals and less ability to pick certain trends. Consider the simulated StO_2 values from Figure 3B. If 40% of the total observations are randomly deleted and the same interaction GAM fitted for the complete dataset is used, the resulting smooths are still able to show a different trend for each group, but because the empirical Bayesian credible intervals for the smooths overlap during the first 3 days with fewer data points, the trend is less pronounced than in the full dataset (Figure 3D). Although the confidence intervals have increased for both smooths, the model still shows different trends with as little as 4 observations per group at certain time points.

6.3 Determination of significance in GAMs for longitudinal data

At the core of a biomedical longitudinal study lies the question of a significant difference between the effect of two or more treatments in different groups. Whereas in rm-ANOVA a *post-hoc* analysis is required to answer such question by calculating some *p-values* after multiple comparisons, GAMs can use a different approach to estimate significance. In essence, the idea behind the estimation of significance in GAMs across different treatment groups is that if the *difference* between the empirical Bayesian confidence intervals of the fitted smooths for such groups is non-zero, then a significant difference exists at that time point(s). The absence of a *p-value* in this case might seem odd, but the empirical Bayesian confidence interval interpretation can be conceptualized in the following manner: Different trends in each group are an indication of an effect by the treatment. This is what happens for the simulated data in Figure 3A, where the chemotherapy causes StO₂ to increase over time.

With this expectation of different trends in each group, computing the difference between the trends will identify if the observed change is significant. The difference between groups with similar trends is likely to yield zero, which would indicate that the treatment is not causing a change in the response in one of the groups (assuming the other group is a Control or Reference group).

Consider the calculation of pairwise differences for the smooths in Figure 3B and Figure 3D. Figure 4 shows the comparison between each treatment group for the full and missing datasets. Here, the “Control” group is used as the reference to which “Treatment” group is being compared. Of notice, the pairwise comparison has been set on the response scale (see Appendix for code details), because otherwise the comparison appears shifted and is not intuitively easy to relate to the original data.

With this correction in mind, the shaded regions over the confidence interval (where it does not cover 0) indicate the time interval where each group has a higher effect than the other. Notice that the shaded region between days 0 and ≈ 2 for the full dataset indicates that through that time, the “Control” group has higher mean StO₂, but as therapy progresses the effect is reversed and by ≈ 3 day it is the “Treatment” group the one that on average, has greater StO₂. This would suggest that the effect of chemotherapy in the “Treatment” group becomes significant after day 3 for the given model. Moreover, notice that although there is no actual measurement at day 3, the model is capable of providing an estimate of when the shift in mean StO₂ occurs.

On the data with missing observations (Figure 3D), the empirical Bayesian credible intervals of the smooths overlap between days 0 and 3. Consequently, the smooth pairwise comparison (Figure 4B) shows that there is no evidence of a significant difference between the groups during that period, but is still able to pick the change on day 3 as the full dataset smooth pairwise comparison.

In a sense, the pairwise smooth comparison is more informative than a *post-hoc p-value*. For biomedical studies, the smooth comparison is able to provide an estimate of *when* and by *how much* a biological process becomes significant. This is advantageous because it can help researchers gain insight on metabolic changes and other biological processes that can be worth examining, and can help refine the experimental design of future studies in order to obtain measurements at time points where a significant change might be expected.

7 Discussion

Biomedical longitudinal non-linear data is particularly challenging to analyze due to the likelihood of missing observations and different correlation structures in the data, which limit the use of rm-ANOVA. Although LMEMs have started to replace rm-ANOVA as the choice to analyze biomedical data, both methods yield biased estimates when they are used to fit non-linear data as we have visually demonstrated in Section 3.5. This “model misspecification” error, also is known as a “Type III” error (Dennis *and others* 2019) is particularly important because although the *p-value* is the common measure of statistical significance, the validity of its interpretation is determined by the agreement of the data and the model. Guidelines for statistical reporting in biomedical journals exist (the SAMPL guidelines) (Lang and Altman 2015) but they have not been widely adopted and in the case of longitudinal data, we consider that researchers would benefit from reporting a visual assessment of the correspondence between the model fit and the data, instead of merely relying on a R^2 value.

In this paper we have presented GAMs as a suitable method to analyze non-linear longitudinal data. It is interesting to note that although GAMs are a well established method to analyze temporal data in different fields (among which are palaeoecology, geochemistry, and ecology) (Hefley *and others* 2017; Pedersen *and others* 2019) they are not routinely used in biomedical research despite an early publication from Hastie and Tibshirani that demonstrated their use in medical research (Hastie and Tibshirani 1995). This is possibly due to the fact that the theory behind GAMs can seem very different from that of rm-ANOVA and LMEMs, but the purpose of Section 4 is to demonstrate that at its core the theory quite simple: Instead of using a linear relationship to model the response (as rm-ANOVA and LMEMs do), GAMs use basis functions to build smooths that are capable of following non-linear trends in the data.

However, from a practical standpoint is equally important to demonstrate how GAMs are computationally implemented. We have provided an example on how GAMs can be fitted using simulated data that follows trends reported in biomedical literature (Vishwanath *and others* 2009) using R and the package *mgcv* (Wood 2017) in Section 6, while a basic workflow for model selection is in the Appendix. One of the features of

GAMs is that their Bayesian interpretation allows to indicate differences between groups without the need of a p -value, and in turn provide a time-based estimate of shifts in the response that can be directly tied to biological values as the pairwise smooth comparisons in Figure 4 indicate. The model is therefore able to provide an estimate of significant change between the groups at time points where data was not directly measured even with missing data exists (\approx day 3 in Figure 4 A, B), which can be used by researchers as feedback on experiment design and to further evaluate important biological changes in future studies.

We have used R as the software of choice for this paper because not only provides a fully developed environment to fit GAMs, but also eases simulation (which is becoming increasingly used for exploratory statistical analysis and power calculations) and provides powerful and convenient methods of visualization, which are key aspects that biomedical researchers might need to consider to make their work reproducible. In this regard, reproducibility is still an issue in biomedical research (Begley and Ioannidis 2015; Weissgerber *and others* 2018), but it is becoming apparent that what other disciplines have experienced in this aspect is likely to impact sooner rather than later this field. Researchers need to plan on how they will make their data, code, and any other materials open and accessible as more journals and funding agencies recognize the importance and benefits of open science in biomedical research. We have made all the data and code used in this paper accessible, and we hope that this will encourage other researchers to do the same with future projects.

8 Conclusion

We have presented GAMs as a method to analyze longitudinal biomedical data. Future directions of this work will include simulation-based estimations of statistical power using GAMs, as well as demonstrating the prediction capabilities of these models using large datasets. By making the data and code used in this paper accessible, we hope to address the need of creating and sharing reproducible work in biomedical research.

9 Acknowledgements

This work was supported by the National Science Foundation Career Award (CBET 1751554, TJM) and the Arkansas Biosciences Institute.

10 Declaration of Conflicting Interests

The Authors declare that there is no conflict of interest.

Supplementary Materials

An Appendix which contains all the code used to create this manuscript, along with a basic workflow to implement GAMs in R is available as Supplementary Material in PDF. A GitHub repository containing all the code used for this paper along with detailed instructions for its use is available at <https://github.com/aimundo/GAMs-biomedical-research>.

11 References

- ABDI, H. (2010). Holm’s sequential Bonferroni procedure. *Encyclopedia of Research Design* **1**, 1–8.
- ALBERS, C. (2019). The problem with unadjusted multiple and sequential statistical testing. *Nature Communications* **10**. doi:10.1038/s41467-019-09941-0.
- BARR, D. J., LEVY, R., SCHEEPERS, C. AND TILY, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language* **68**, 255–278.
- BECK, N. AND JACKMAN, S. (1998). Beyond linearity by default: Generalized additive models. *American Journal of Political Science* **42**, 596.
- BEGLEY, C. G. AND IOANNIDIS, J. P. A. (2015). Reproducibility in science. *Circulation Research* **116**, 116–126.
- CHARAN, J. AND KANTHARIA, N. (2013). How to calculate sample size in animal studies? *Journal of Pharmacology and Pharmacotherapeutics* **4**, 303.
- DEMIDOV, V., MAEDA, A., SUGITA, M., MADGE, V., SADANAND, S., FLUERARU, C. AND VITKIN, I. A. (2018). Preclinical longitudinal imaging of tumor microvascular radiobiological response with functional optical coherence tomography. *Scientific Reports* **8**. doi:10.1038/s41598-017-18635-w.
- DENNIS, B., PONCIANO, J. M., TAPER, M. L. AND LELE, S. R. (2019). Errors in statistical inference under model misspecification: Evidence, hypothesis testing, and AIC. *Frontiers in Ecology and Evolution* **7**. doi:10.3389/fevo.2019.00372.
- GEISSER, S. AND GREENHOUSE, S. W. (1958). An extension of box’s results on the use of the f distribution in multivariate analysis. *The Annals of Mathematical Statistics* **29**, 885–891.
- GELMAN, A., HILL, J. AND YAJIMA, M. (2012). Why we (usually) don’t have to worry about multiple comparisons. *Journal of Research on Educational Effectiveness* **5**, 189–211.
- GREENHOUSE, S. W. AND GEISSER, S. (1959). On methods in the analysis of profile data. *Psychometrika* **24**, 95–112.
- GREENING, G. J., MILLER, K. P., SPAINHOUR, C. R., CATO, M. D. AND MULDOON, T. J. (2018). Effects of isoflurane anesthesia on physiological parameters in murine subcutaneous tumor allografts measured via diffuse reflectance spectroscopy. *Biomedical Optics Express* **9**, 2871.
- GUEORGUEVA, R. AND KRYSTAL, J. H. (2004). Move over ANOVA. *Archives of General Psychiatry* **61**,

- HALSEY, L. G., CURRAN-EVERETT, D., VOWLER, S. L. AND DRUMMOND, G. B. (2015). The fickle p value generates irreproducible results. *Nature Methods* **12**, 179–185.
- HAREZLAK, J., RUPPERT, D. AND WAND, M. P. (2018). Semiparametric regression with r. Springer New York. doi:10.1007/978-1-4939-8853-2.
- HASTIE, T. AND TIBSHIRANI, R. (1995). Generalized additive models for medical research. *Statistical Methods in Medical Research* **4**, 187–196.
- HASTIE, T. AND TIBSHIRANI, R. (1987). Generalized additive models: Some applications. *Journal of the American Statistical Association* **82**, 371–386.
- HAVERKAMP, N. AND BEAUDUCEL, A. (2017). Violation of the sphericity assumption and its effect on type-i error rates in repeated measures ANOVA and multi-level linear models (MLM). *Frontiers in Psychology* **8**. doi:10.3389/fpsyg.2017.01841.
- HEFLEY, T. J., BROMS, K. M., BROST, B. M., BUDERMAN, F. E., KAY, S. L., SCHARF, H. R., TIPTON, J. R., WILLIAMS, P. J. AND HOOTEN, M. B. (2017). The basis function approach for modeling autocorrelation in ecological data. *Ecology* **98**, 632–646.
- HUYNH, H. AND FELDT, L. S. (1976). Estimation of the box correction for degrees of freedom from sample data in randomized block and split-plot designs. *Journal of Educational Statistics* **1**, 69–82.
- JONES, J. D., RAMSER, H. E., WOESSNER, A. E. AND QUINN, K. P. (2018). In vivo multiphoton microscopy detects longitudinal metabolic changes associated with delayed skin wound healing. *Communications Biology* **1**. doi:10.1038/s42003-018-0206-4.
- KAMSTRA, J. I., DIJKSTRA, P. U., LEEUWEN, M. VAN, ROODENBURG, J. L. N. AND LANGENDIJK, J. A. (2015). Mouth opening in patients irradiated for head and neck cancer: A prospective repeated measures study. *Oral Oncology* **51**, 548–555.
- KESELMAN, H. J., ALGINA, J. AND KOWALCHUK, R. K. (2001). The analysis of repeated measures designs: A review. *British Journal of Mathematical and Statistical Psychology* **54**, 1–20.
- LAIRD, N. M. AND WARE, J. H. (1982). Random-effects models for longitudinal data. *Biometrics* **38**, 963.
- LANG, T. A. AND ALTMAN, D. G. (2015). Basic statistical reporting for articles published in biomedical journals: The “statistical analyses and methods in the published literature” or the SAMPL guidelines. *International Journal of Nursing Studies* **52**, 5–9.

LIU, C., CRIPE, T. P. AND KIM, M.-O. (2010). Statistical issues in longitudinal data analysis for treatment efficacy studies in the biomedical sciences. *Molecular Therapy* **18**, 1724–1730.

MA, Y., MAZUMDAR, M. AND MEMTSOUDIS, S. G. (2012). Beyond repeated-measures analysis of variance. *Regional Anesthesia and Pain Medicine* **37**, 99–105.

MARRA, G. AND WOOD, S. N. (2012). Coverage properties of confidence intervals for generalized additive model components. *Scandinavian Journal of Statistics* **39**, 53–74.

MAXWELL, S. E., DELANEY, H. D. AND KELLEY, K. (2017). Designing experiments and analyzing data. Routledge. doi:10.4324/9781315642956.

MCELREATH, R. (2018). Statistical rethinking. Chapman; Hall/CRC. doi:10.1201/9781315372495.

MILLER, D. L. (2019). Bayesian views of generalized additive modelling. *arXiv preprint arXiv:1902.01330*.

MOLENBERGHS, G. (2004). Analyzing incomplete longitudinal clinical trial data. *Biostatistics* **5**, 445–464.

NAKAGAWA, S. (2004). A farewell to bonferroni: The problems of low statistical power and publication bias. *Behavioral Ecology* **15**, 1044–1045.

NELDER, J. A. AND WEDDERBURN, R. W. M. (1972). Generalized linear models. *Journal of the Royal Statistical Society. Series A (General)* **135**, 370.

PAVLOV, M. V., KALGANOVA, T. I. AND LYUBIMTSEVA, Y. S. (2018). Multimodal approach in assessment of the response of breast cancer to neoadjuvant chemotherapy. *Journal of Biomedical Optics* **23**, 1.

PEDERSEN, E. J., MILLER, D. L., SIMPSON, G. L. AND ROSS, N. (2019). Hierarchical generalized additive models in ecology: An introduction with mgcv. *PeerJ* **7**, e6876.

PINHEIRO, J. AND BATES, D. (2006). Mixed-effects models in S and S-PLUS. Springer Science; Business Media. doi:https://doi.org/10.1007/b98882.

PINHEIRO, J., BATES, D., DEBROY, S., SARKAR, D. AND R CORE TEAM. (2020). nlme: Linear and nonlinear mixed effects models. <https://CRAN.R-project.org/package=nlme>.

POTTHOFF, R. F., TUDOR, G. E., PIEPER, K. S. AND HASSELBLAD, V. (2006). Can one assess whether missing data are missing at random in medical studies? *Statistical Methods in Medical Research* **15**, 213–234.

R CORE TEAM. (2020). R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.

- 616 RITTER, G., COHEN, L., WILLIAMS, C., RICHARDS, E., OLD, L. AND WELT, S. (2001). Serological
617 analysis of human anti-human antibody responses in colon cancer patients treated with repeated doses of
618 humanized monoclonal antibody A33. *Cancer Research* **61**, 6851–6859.
- 619 ROBLYER, D., UEDA, S., CERUSSI, A., TANAMAI, W., DURKIN, A., MEHTA, R., HSIANG, D., BUTLER, J.
620 A., MCLAREN, C., CHEN, W.-P., ET AL. (2011). Optical imaging of breast cancer oxyhemoglobin flare
621 correlates with neoadjuvant chemotherapy response one day after starting treatment. *Proceedings of the*
622 *National Academy of Sciences* **108**, 14626–14631.
- 623 ROSE, N. L., YANG, H., TURNER, S. D. AND SIMPSON, G. L. (2012). An assessment of the mechanisms
624 for the transfer of lead and mercury from atmospherically contaminated organic soils to lake sediments
625 with particular reference to scotland, UK. *Geochimica et Cosmochimica Acta* **82**, 113–135.
- 626 ROTH, E. M., GOLDBERG, A. C., CATAPANO, A. L., TORRI, A., YANCOPOULOS, G. D., STAHL, N.,
627 BRUNET, A., LECORPS, G. AND COLHOUN, H. M. (2017). Antidrug antibodies in patients treated with
628 alirocumab. *New England Journal of Medicine* **376**, 1589–1590.
- 629 SCHEFFER, J. (2002). Dealing with missing data. *Research Letters in the Information and Mathematical*
630 *Sciences* **3**, 153–160.
- 631 SCHOBER, P. AND VETTER, T. R. (2018). Repeated measures designs and analysis of longitudinal data.
632 *Anesthesia & Analgesia* **127**, 569–575.
- 633 SIMPSON, G. L. (2020). Gratia: Graceful 'ggplot'-based graphics and other functions for GAMs fitted using
634 'mgcv'. <https://CRAN.R-project.org/package=gratia>.
- 635 SIMPSON, G. L. (2018). Modelling palaeoecological time series using generalised additive models. *Frontiers*
636 *in Ecology and Evolution* **6**. doi:10.3389/fevo.2018.00149.
- 637 SIO, T. T., ATHERTON, P. J., BIRCKHEAD, B. J., SCHWARTZ, D. J., SLOAN, J. A., SEISLER, D. K.,
638 MARTENSON, J. A., LOPRINZI, C. L., GRIFFIN, P. C., MORTON, R. F., ET AL. (2016). Repeated
639 measures analyses of dermatitis symptom evolution in breast cancer patients receiving radiotherapy in a
640 phase 3 randomized trial of mometasone furoate vs placebo (N06C4 [alliance]). *Supportive Care in Cancer*
641 **24**, 3847–3855.
- 642 SKALA, M. C., FONTANELLA, A., LAN, L., IZATT, J. A. AND DEWHIRST, M. W. (2010). Longitudinal
643 optical imaging of tumor metabolism and hemodynamics. *Journal of Biomedical Optics* **15**, 011112.
- 644 TANK, A., PETERSON, H. M., PERA, V., TABASSUM, S., LEPROUX, A., O'SULLIVAN, T., JONES, E.,
645 CABRAL, H., KO, N., MEHTA, R. S., ET AL. (2020). Diffuse optical spectroscopic imaging reveals

distinct early breast tumor hemodynamic responses to metronomic and maximum tolerated dose regimens.

Breast Cancer Research **22**. doi:10.1186/s13058-020-01262-1.

UGRINOWITSCH, C., FELLINGHAM, G. W. AND RICARD, M. D. (2004). Limitations of ordinary least squares models in analyzing repeated measures data. *Medicine & Science in Sports & Exercise*, 2144–2148.

VISHWANATH, K., YUAN, H., BARRY, W. T., DEWHIRST, M. W. AND RAMANUJAM, N. (2009). Using optical spectroscopy to longitudinally monitor physiological changes within solid tumors. *Neoplasia* **11**, 889–900.

WAGENMAKERS, E.-J., LEE, M., LODEWYCKX, T. AND IVERSON, G. J. (2008). Bayesian versus frequentist inference. In Bayesian evaluation of informative hypotheses. Springer New York. pp 181–207.

WANG, B., ZHOU, Z., WANG, H., TU, X. M. AND FENG, C. (2019). The p-value and model specification in statistics. *General Psychiatry* **32**, e100081.

WEGMAN, E. J. AND WRIGHT, I. W. (1983). Splines in statistics. *Journal of the American Statistical Association* **78**, 351–365.

WEISS, R. E. (2005). Modeling longitudinal data. Springer New York. doi:10.1007/0-387-28314-5.

WEISSGERBER, T. L., GARCIA-VALENCIA, O., GAROVIC, V. D., MILIC, N. M. AND WINHAM, S. J. (2018). Why we need to report more than data were analyzed by t-tests or ANOVA. *eLife* **7**. doi:10.7554/elife.36163.

WEST, B. T., WELCH, K. B. AND GALECKI, A. T. (2014). Linear mixed models: A practical guide using statistical software, second edition. Taylor & Francis. <https://books.google.com/books?id=hjT6AwAAQBAJ>.

WOLFINGER, R. D. (1996). Heterogeneous variance: Covariance structures for repeated measures. *Journal of Agricultural, Biological, and Environmental Statistics* **1**, 205.

WOOD, S. N. (2017). Generalized additive models. Chapman; Hall/CRC. doi:10.1201/9781315370279.

WOOD, S. N., PYA, N. AND SÄFKEN, B. (2016). Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association* **111**, 1548–1563.

YANG, L., QIN, G., ZHAO, N., WANG, C. AND SONG, G. (2012). Using a generalized additive model with autoregressive terms to study the effects of daily temperature on mortality. *BMC Medical Research Methodology* **12**. doi:10.1186/1471-2288-12-165.

List of Figures

- 1 Simulated responses from two groups with correlated errors using a LMEM and a rm-ANOVA model. Top row: linear response, bottom row: quadratic response. A: Simulated linear data with known mean response (thin lines) and individual responses (points) showing the dispersion of the data. D: Simulated quadratic data with known mean response (thin lines) and individual responses (points) showing the dispersion of the data. B,E: Estimates from the rm-ANOVA model for the mean group response (linear or quadratic). Points represent the original raw data. The rm-ANOVA model not only fails to pick the trend of the quadratic data (D) but also assigns a global estimate that does not take between-subject variation. C, F: Estimates from the LMEM in the linear and quadratic case. The LMEM incorporates a random effect for each subject, but this model and the rm-ANOVA model are unable to follow the trend of the data and grossly bias the initial estimates for each group in the quadratic case (bottom row). 26
- 2 Basis functions for a single smoother for time with five knots. A: Basis functions for a single smoother for time for the simulated data of Group 1 from Figure 2. B: Matrix of basis function weights. Each basis function is multiplied by a coefficient which can be positive or negative. The coefficient determines the overall effect of each basis in the final smoother. C: Weighted basis functions. Each of the four basis functions of panel A has been weighted by the corresponding coefficient shown in Panel B. Note the corresponding increase (or decrease) in magnitude of each weighted basis function. D: Smoother for time and original data points. The smoother (line) is the result of the sum of each weighted basis function at each time point, with simulated values for the group shown as points. 27
- 3 Simulated data and smooths for oxygen saturation in tumors. A: Simulated data that follows previously reported trends (inset) in tumors under chemotherapy (Treatment) or saline (Control) treatment. Simulated data is from a normal distribution with standard deviation of 10% with 10 observations per time point. Lines indicate mean oxygen saturation B: Smooths from the GAM model for the full simulated data with interaction of Group and Treatment. Lines represent trends for each group, shaded regions are 95% confidence intervals. C: The rm-ANOVA model for the simulated data, which does not capture the changes in each group over time. D: Smooths for the GAM model for the simulated data with 40% of its observations missing. Lines represent trends for each group, shaded regions are 95% empirical Bayesian confidence intervals. 28
- 4 Pairwise comparisons for smooth terms. A: Pairwise comparisons for the full dataset. B: Pairwise comparisons for the dataset with missing observations. Significant differences exist where the 95% empirical Bayesian credible interval does not cover 0. In both cases the effect of treatment is significant after day 3. 29

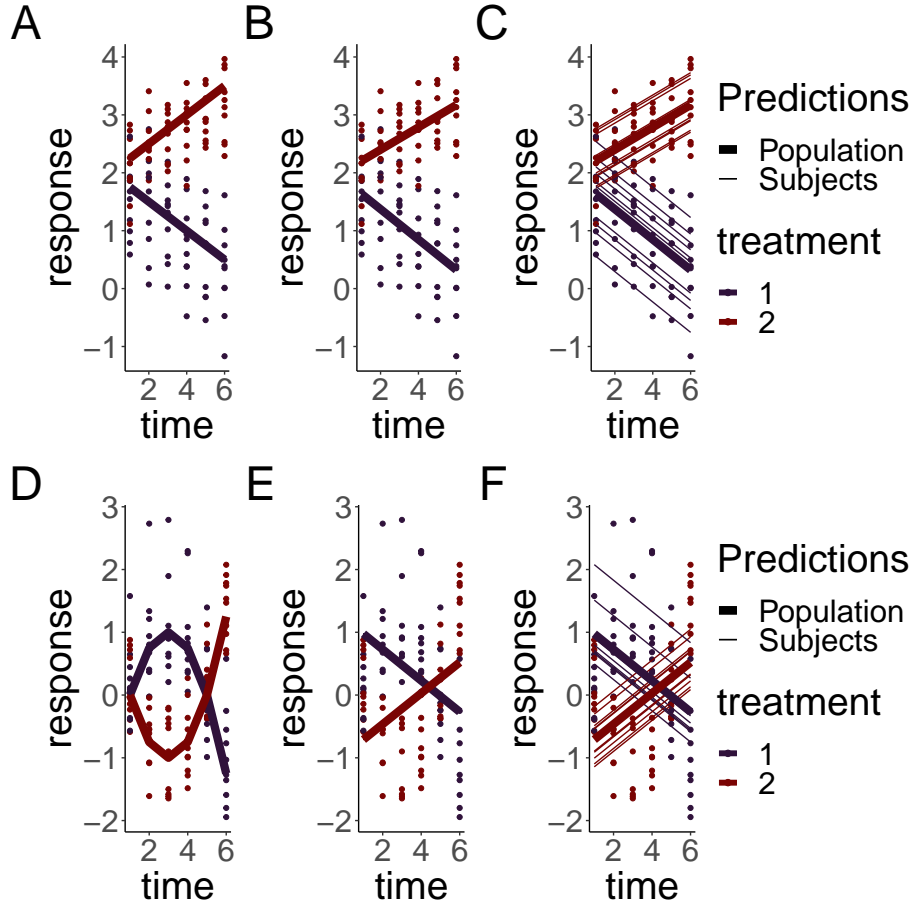


Figure 1: Simulated responses from two groups with correlated errors using a LMEM and a rm-ANOVA model. Top row: linear response, bottom row: quadratic response. A: Simulated linear data with known mean response (thin lines) and individual responses (points) showing the dispersion of the data. D: Simulated quadratic data with known mean response (thin lines) and individual responses (points) showing the dispersion of the data. B,E: Estimates from the rm-ANOVA model for the mean group response (linear of quadratic). Points represent the original raw data. The rm-ANOVA model not only fails to pick the trend of the quadratic data (D) but also assigns a global estimate that does not take between-subject variation. C, F: Estimates from the LMEM in the linear and quadratic case. The LMEM incorporates a random effect for each subject, but this model and the rm-ANOVA model are unable to follow the trend of the data and grossly bias the initial estimates for each group in the quadratic case (bottom row).

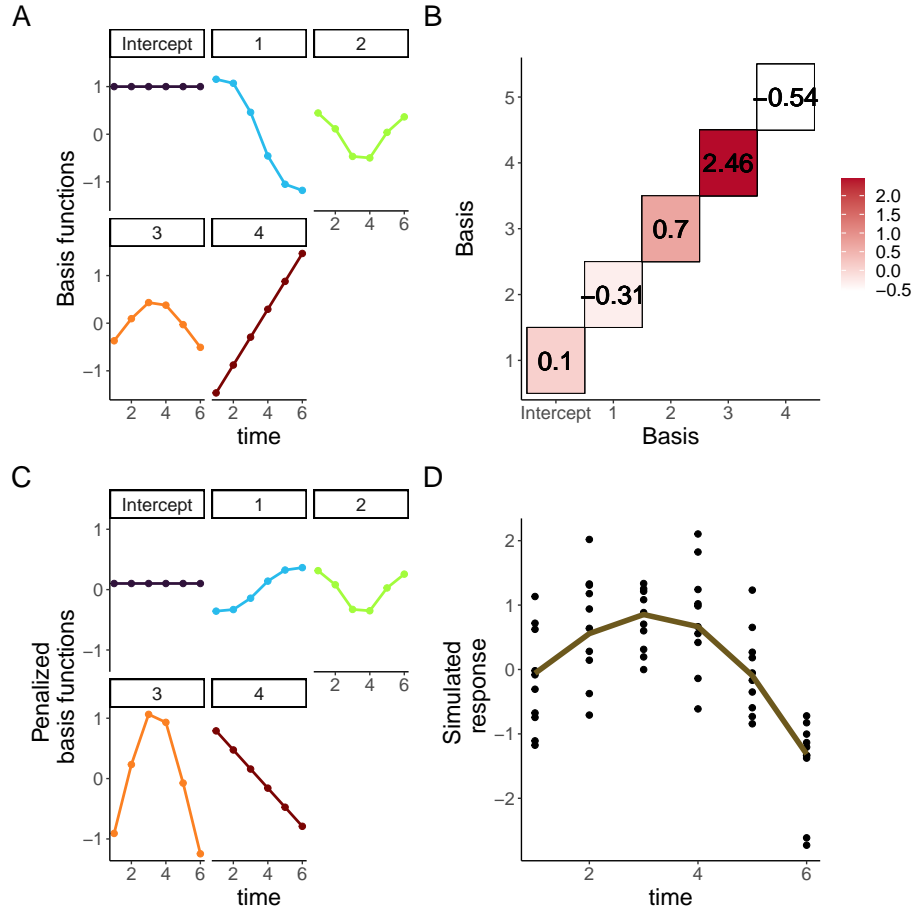


Figure 2: Basis functions for a single smoother for time with five knots. A: Basis functions for a single smoother for time for the simulated data of Group 1 from Figure 2. B: Matrix of basis function weights. Each basis function is multiplied by a coefficient which can be positive or negative. The coefficient determines the overall effect of each basis in the final smoother. C: Weighted basis functions. Each of the four basis functions of panel A has been weighted by the corresponding coefficient shown in Panel B. Note the corresponding increase (or decrease) in magnitude of each weighted basis function. D: Smoother for time and original data points. The smoother (line) is the result of the sum of each weighted basis function at each time point, with simulated values for the group shown as points.

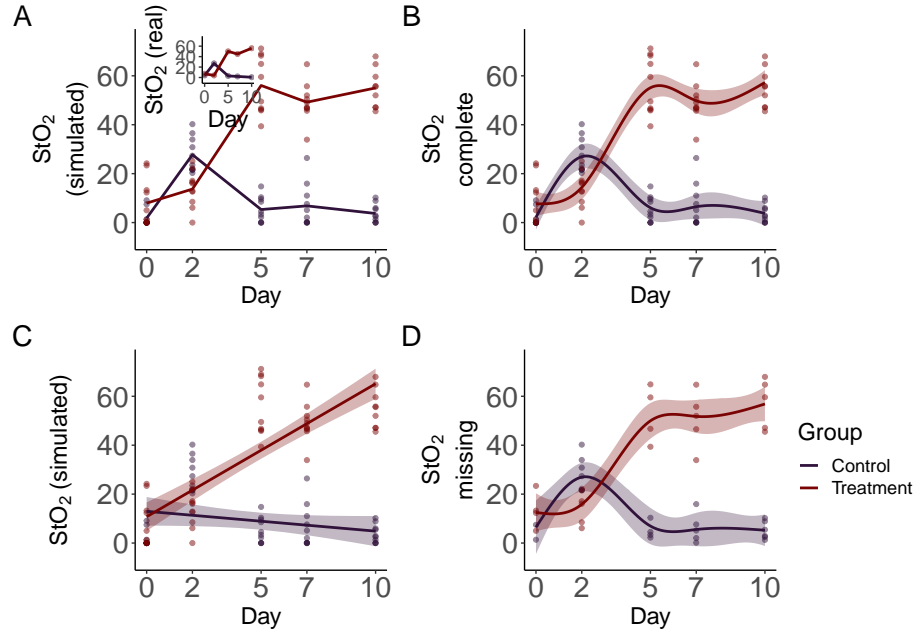


Figure 3: Simulated data and smooths for oxygen saturation in tumors. A: Simulated data that follows previously reported trends (inset) in tumors under chemotherapy (Treatment) or saline (Control) treatment. Simulated data is from a normal distribution with standard deviation of 10% with 10 observations per time point. Lines indicate mean oxygen saturation B: Smooths from the GAM model for the full simulated data with interaction of Group and Treatment. Lines represent trends for each group, shaded regions are 95% confidence intervals. C: The rm-ANOVA model for the simulated data, which does not capture the changes in each group over time. D: Smooths for the GAM model for the simulated data with 40% of its observations missing. Lines represent trends for each group, shaded regions are 95% empirical Bayesian confidence intervals.

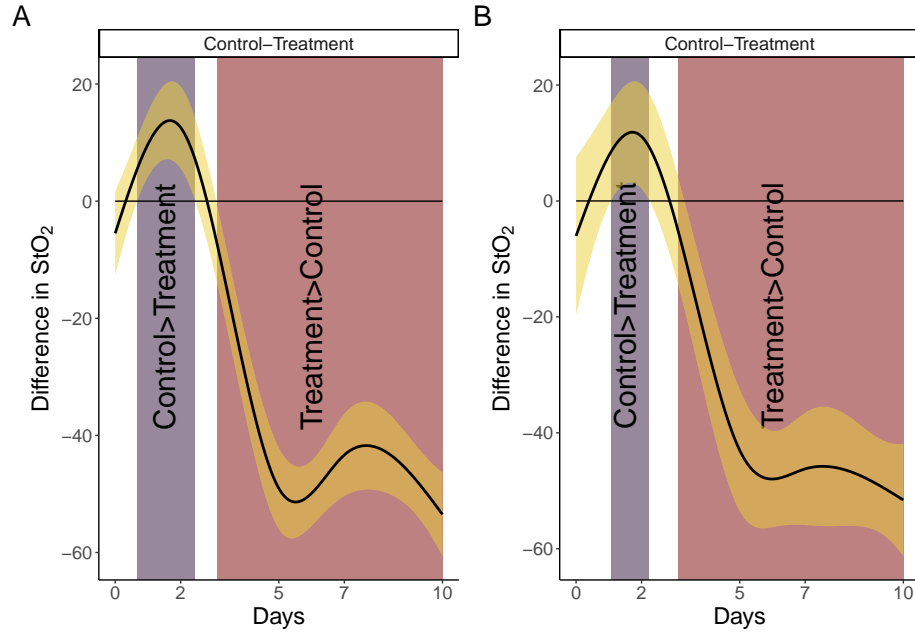


Figure 4: Pairwise comparisons for smooth terms. A: Pairwise comparisons for the full dataset. B: Pairwise comparisons for the dataset with missing observations. Significant differences exist where the 95% empirical Bayesian credible interval does not cover 0. In both cases the effect of treatment is significant after day 3.