

The statistical analysis of non-linear longitudinal data in biomedical research using generalized additive models

Beyond repeated measures ANOVA and Linear Mixed Models

Ariel Mundo* Timothy J. Muldoon† John R. Tipton‡

Contents

1	Background	1
2	Challenges presented by longitudinal studies	4
2.1	The repeated measures ANOVA	4
2.2	Linear relationship	4
2.3	Covariance in rm-ANOVA and LMEMs	5
2.4	Missing observations	6
2.5	How does an rm-ANOVA fit looks like? A visual representation using simulation	6
3	GAMs as a special case of Generalized Linear Models	8
3.1	GAMs and Basis Functions	8
3.2	GAMs and covariance	10
3.3	Determination of significance in GAMs	10
4	References	11
A	Simulation	14
A.1	Compound symmetry and independent errors in linear and quadratic responses	14

1 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

*Department of Biomedical Engineering, University of Arkansas, Fayetteville

†Department of Biomedical Engineering, University of Arkansas, Fayetteville

‡Department of Mathematical Sciences, University of Arkansas, Fayetteville

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 [1–4], antibody expression [5,6], and cell metabolism [7,8] 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[9], to weekly measurements to analyze a mid-term response like the evolution of dermatitis symptoms in breast cancer patients [10], to monthly measurements to study a long-term response such as mouth opening following RT in neck cancer patients [11].

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 a limiting frequency [12] 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*) [13,14].

The expected linear behavior of the response through time is a key requisite in rm-ANOVA [15]. 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 in preclinical and clinical settings [1,8,16]. 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 highly-variable data the estimates are inevitably biased, because the model is only able to accommodate linear trends that are far from adequately representing the biological phenomenon of interest.

A *post hoc* analysis is the statistical test 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[17]. 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 [18].

Additionally, the *p-value* itself is highly variable, and multiple comparisons can inflate the false positivity rate (Type I error or α) [19,20], consequently biasing the conclusions of the study. Corrections exist to address the Type I error issue of multiple comparisons (such as Bonferroni [21]), but they in turn reduce statistical power ($1 - \beta$)[22], and lead to increased Type II error (failing to reject the null hypothesis when the null hypothesis is false) [23,24]. Therefore, the tradeoff of *post hoc* comparisons in rm-ANOVA between Type I, II and III errors might be difficult to balance 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 [25]. Corrections can be made in rm-ANOVA in the absence of compound symmetry [26,27], but the effectiveness the correction is limited by the size of the sample, the number of measurements[28], and group sizes [29]. In the case of biomedical research, where living subjects are frequently used, sample sizes are often not “large” due to ethical and budgetary reasons [30] 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 [13]. 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 information. 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 [8,16]. 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 [15]. However, LMEMs impose restrictions in the distribution of the errors of the random effects, which need to be normally distributed and independent [13,31]. And even more importantly, LMEMs also expect a linear relationship between the response and time [15], 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 non-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 [32–34], health-environment interactions [35] and the dynamics of government in political science [36]. 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 [37] 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[38] has an extensive collection of documentation and functions to fit GAMs in the package *mgcv* [37,39] 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 [28]. 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 is 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, using simulated data that reproduces patterns in previously reported studies [16] we illustrate the type of non-linear longitudinal data that often occurs in biomedical research. 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 by improving the standards for reproducibility in biomedical research.

2 Challenges presented by longitudinal studies

2.1 The repeated measures ANOVA

The *repeated measures analysis of variance* (rm-ANOVA) is the standard statistical analysis for longitudinal data in biomedical research. This statistical methodology requires 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.

2.2 Linear relationship

2.2.1 The repeated measures ANOVA case

In a biomedical longitudinal study, two or more groups of subjects (e.g., patients, mice, samples) are subject to different treatments (e.g., a “treatment” group receives a novel drug 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 in the “treatment” group.

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

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

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_{tij} represent random variation not explained by the *fixed* effects, and are assumed to be $\sim N(0, \sigma^2)$. 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). The linear models then can be expressed as

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

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

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

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

2.2.2 The Linear Mixed Model Case

A linear mixed model (LMEM) is a class of statistical model that incorporates *fixed effects* to model the relationship between the covariates and the response, and *random effects* to model subject variability that is not the primary focus of the study but that might be important to distinguish [15,40]. A LMEM with interaction between time and treatment for a longitudinal study can be written as:

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

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

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 independently normally distributed with mean zero and variance σ_μ^2 , which can be expressed as $\mu_{ij} \sim N(0, \sigma_\mu^2)$. In essence, the *random effect* in a LMEM enables to fit models with different slopes at the subject-level [15]. 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.

2.3 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 [41] (For an in-depth analysis of the covariance matrix see [40,42]).

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 [26,42,43]. However, it has been shown that this condition is frequently not justified because the correlation between measurements tends to change over time [44]; and it is higher between consecutive measurements [13,25]. Although corrections can be made (such as Huynh-Feldt or Greenhouse-Geisser) [26,27] the effectiveness of each correction is limited because it depends on the size of the sample, the number of repeated measurements [28], and they are not robust if the group sizes are unbalanced [29]. 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 [15]. Nevertheless, the analysis required to determine an appropriate variance-covariance structure for the data can be a long 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.

2.4 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 [45]. 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) [42]. 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 [46]. If the data are MCAR, this means that the missingness is completely unrelated to the collected information [47], 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 [48].

In the case of LMEMs, inferences from the model are valid when missing observations in the data exist that are MAR or MCAR [40]. 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.

This section has presented the assumptions for analyzing longitudinal data using rm-ANOVA and LMEMs and compared their differences regarding linearity, the covariance matrix and missing data. Of notice, LMEMs offer a more robust and flexible approach than rm-ANOVA and if the data follows a linear trend, they provide an excellent choice to derive inferences from a repeated measures study. However, when the data presents high a non-linear behavior, LMEMs and rm-ANOVA fail to capture the trend of the data. To better convey the issues of linearity and correlation in linear models fitted to non-linear data, simulation is used in the next section.

2.5 How does an rm-ANOVA fit looks like? A visual representation using simulation

To 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)) is fitted to each group, using R[38] and the package *nlme*[49]. 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 with different intercepts and slopes; a negative slope is used for Group 1 and a positive slope is used for Group 2 (Figure @ref(fig:l-q-response, A)). 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 1, C). 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.

Specifically, the rationale for the chosen linear and quadratic functions is the likelihood that a measured response in two treatment groups is similar in the initial phase of the study, but as treatment progresses a divergence in the trend of the response indicates a difference in the effect of each treatment. 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 for linear and non-linear longitudinal data is in Figure 3 and Figure 4 in the Appendix, 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.

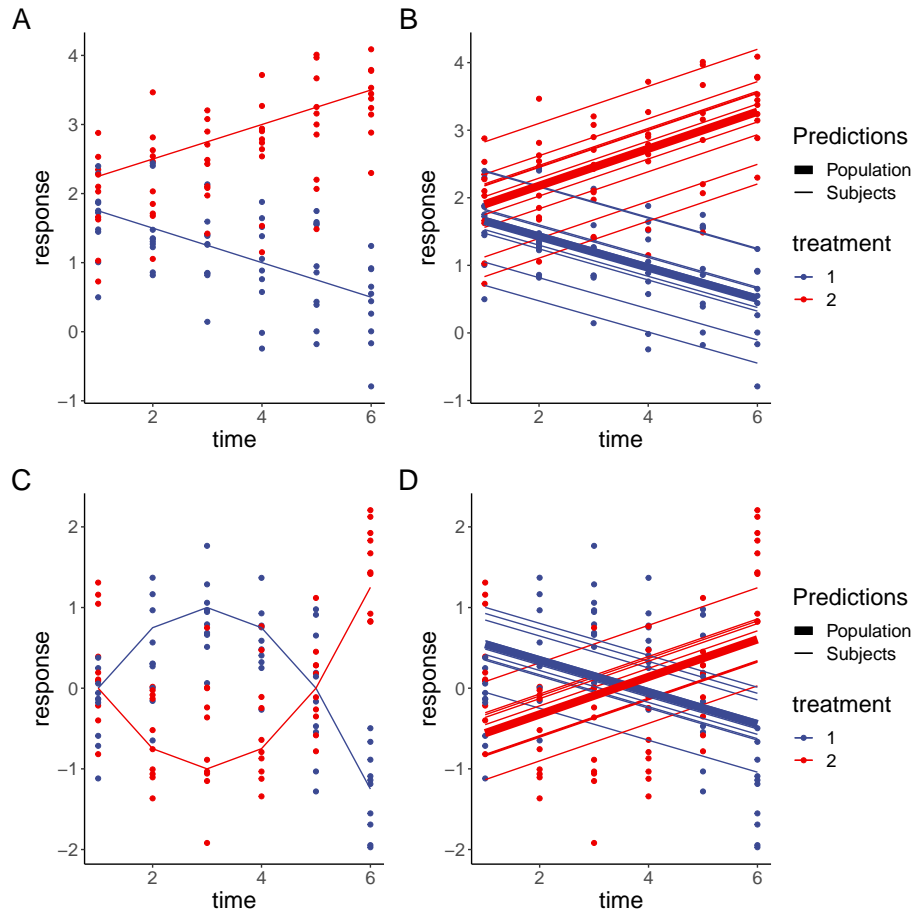


Figure 1: Simulated linear responses from two groups with correlated (top row) or independent (bottom row) errors using a rm-ANOVA model. A, C: Simulated data with known mean response (linear or quadratic, thin lines) and individual responses (points) showing the dispersion of the data. B, D: Estimations from the rm-ANOVA model for the mean group response (linear or quadratic). Thick lines are the predicted mean response per group, thin lines are the random effects for each subject and points represent the original raw data. The rm-ANOVA model does not pick the trend of the quadratic data.

The simulation shows that the fit produced by the rm-ANOVA model is good for linear data, as the predictions for the mean response are reasonably close to the “truth” of the simulated data (Figure 1, B). When the linearity and compound symmetry assumptions are met, the model approximates well the individual trends and the mean trends by group.

However, consider the case when the data follows a non-linear trend, such as the simulated data in Figure 1, C. Here, the mean response per group was simulated using a quadratic function but errors, individual responses and the rm-ANOVA model were produced in the same manner as in 1 A and B. The mean response in the simulated data with quadratic behavior is changing 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

(1) to this data produces the fit that appears in panel D in Figure 1.

A comparison of the fitted mean response of the rm-ANOVA model to the simulated data in Figure (1, D indicates that the model is not capturing the changes within each group in a good way. Specifically, note that the fitted mean response of the rm-ANOVA model (panel D) 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 by the model. Moreover, the rm-ANOVA model is not being able to capture the fact that the initial values are the same in each group, and instead fits non-parallel lines that have initial values that are markedly different from the “true” initial values in each case (compare panels C and D). If such change has important physiological implications, the rm-ANOVA model omits it from the fitted mean response, potentially 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. Although the model fitted to the simulated data was an rm_ANOVA model, the main issue of an expected linear trend in the response is the same in the case of a LMEM. In the following section, we present generalized additive models (GAMs) as an alternative to analyze longitudinal non-linear data.

3 GAMs as a special case of Generalized Linear Models

3.1 GAMs and Basis Functions

Generalized linear models (GLMs) are a family of models that fit a linear response function to data that do not have normally distributed errors[50]. 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[34,37,51]. 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 function $f(x_t | \beta_j)$ 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 will estimate a linear relationship if the data is consistent with a linear response. One possible function for $f(x_t | \beta_j)$ that allows for non-linear responses is a polynomial, 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 the fit [36]. 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 parametric relationship that can be fit within the GLM framework, by using *basis functions* expansions of the covariates and by estimating random coefficients for these basis functions and this step is achieved by using *basis functions* to represent them. A *basis* is a set of functions that spans the space where the smooths that approximate $f(x_t | \beta_j)$ exist [34]. 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 [52], which in the case of Equation (1) is restricted to a linear family of functions. In the case of Equation (5), the basis function is $f(x_t | \beta_j)$, which means that the model allows relationships beyond linear for the covariates.

A commonly used *basis function* is the cubic spline, which is a smooth curve constructed from cubic polynomials joined together at the knot locations[34,37]. Cubic splines have a long history in solving non-parametric

statistical problems and are often a default choice to fit GAMs as they are a simple, flexible and powerful option to obtain visual smoothness [53]. 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 @fig:l-q-response), C. 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 the timeline. The timeline can be divided equally spaced *knots*, each knot being a region where a different basis function 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 already optimized in the package *mgcv* depending on the number of knots. In Panel A of Figure 2, 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 wiggliness of the fit, each of the basis functions of Panel A is penalized by multiplying it by a coefficient according to the penalty matrix of Panel B.

In other words, the six points of each basis are multiplied by the corresponding coefficient in panel B, thereby increasing or decreasing the original basis functions of Panel A. In Figure 2, Panel C shows the resulting penalized basis functions. Note that the penalization for basis 1 has resulted in a decrease of its overall value (because the coefficient for that basis function is negative and less than 1); on the other hand, basis 3 has roughly doubled its value. Finally, the penalized basis functions are added at each timepoint to produce the smooth term. The resulting smooth term for the effect of *time* is shown in Panel D (orange line) along the simulated values per group, which appear as points.

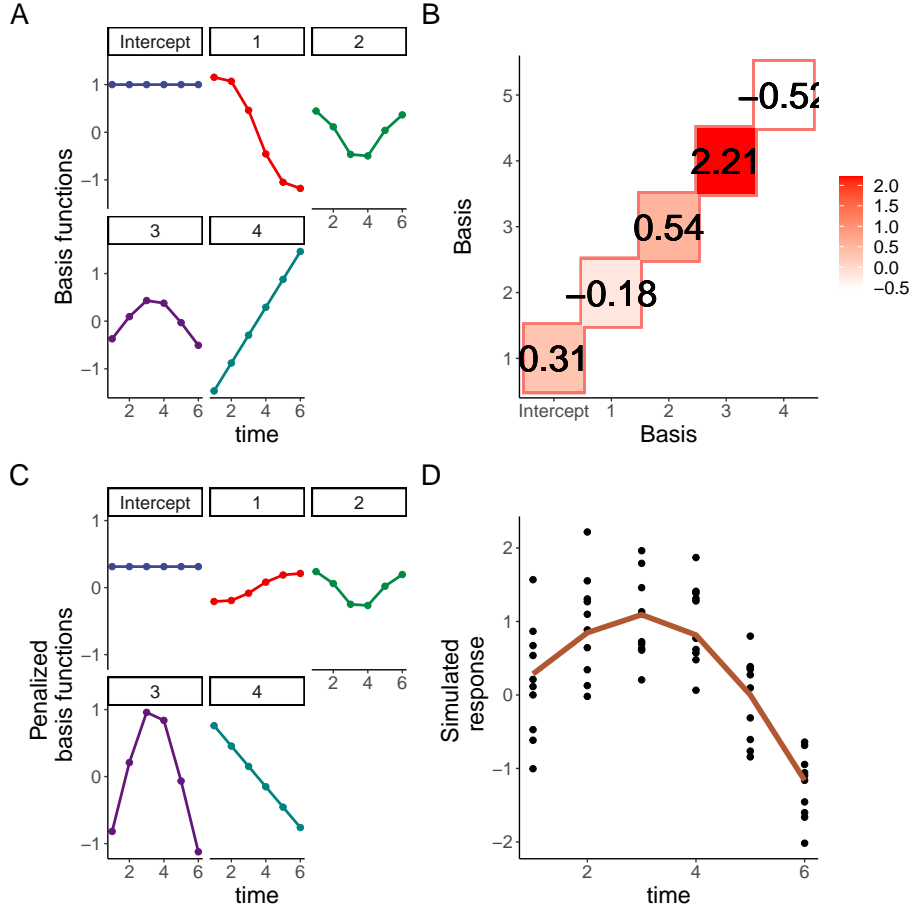


Figure 2: Basis functions for a single smoother for time with five knots. A: Basis functions for a single smoother for time for the simulated data of Group 1 from Figure 2, the intercept basis is not shown. B: Penalty matrix for the basis functions. Each basis function is penalized by a coefficient which can be positive or negative. The coefficient determines the overall effect of each basis in the final smoother. C: Penalized basis functions. Each of the four basis functions of panel A has been penalized by the corresponding coefficient shown in Panel B, note the corresponding increase (or decrease) of each basis. D: Smoother for time and original datapoints. The smoother (line) is the result of the sum of each penalized basis function at each time point, simulated values for the group appear as points.

3.2 GAMs and covariance

Although the specific methods of how GAMs model correlation structures is a topic beyond the scope of this work, it suffices to say that GAMs can handle correlation structures beyond compound symmetry. A detailed description on basis functions and correlations can be found in [52]. In a practical sense, when a GAM is implemented for longitudinal data, a spline can be added to the model for the *time* effect to account for the repeated measures over time. An example where this is covered is in the Appendix.

3.3 Determination of significance in GAMs

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 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 confidence intervals of the fitted smooths for such

groups is non-zero, then a significant difference exists at that timepoint (or timepoints). The absence of a *p-value* in this case might seem odd, but when a confidence interval comparison is put into context the validity of its rationale becomes apparent. The major advantage of the significance estimation in GAMs is that they allow to determine significance at a specific timepoint, rather than providing a single *p-value* for the overall treatment effect. The estimation of significant differences in GAMs is covered in detail in the Appendix using simulated longitudinal data that follows previously reported trends of tumor response between multiple treatment groups [16].

4 References

- [1] D. Roblyer, S. Ueda, A. Cerussi, W. Tanamai, A. Durkin, R. Mehta, D. Hsiang, J.A. Butler, C. McLaren, W.-P. Chen, others, Optical imaging of breast cancer oxyhemoglobin flare correlates with neoadjuvant chemotherapy response one day after starting treatment, *Proceedings of the National Academy of Sciences*. 108 (2011) 14626–14631.
- [2] A. Tank, H.M. Peterson, V. Pera, S. Tabassum, A. Leproux, T. O’Sullivan, E. Jones, H. Cabral, N. Ko, R.S. Mehta, others, Diffuse optical spectroscopic imaging reveals distinct early breast tumor hemodynamic responses to metronomic and maximum tolerated dose regimens, *Breast Cancer Research*. 22 (2020) 1–10.
- [3] M.V. Pavlov, T.I. Kalganova, Y.S. Lyubimtseva, V.I. Plekhanov, G.Y. Golubyatnikov, O.Y. Ilyinskaya, A.G. Orlova, P.V. Subochev, D.V. Safonov, N.M. Shakhova, others, Multimodal approach in assessment of the response of breast cancer to neoadjuvant chemotherapy, *Journal of Biomedical Optics*. 23 (2018) 091410.
- [4] V. Demidov, A. Maeda, M. Sugita, V. Madge, S. Sadanand, C. Flueraru, I.A. Vitkin, Preclinical longitudinal imaging of tumor microvascular radiobiological response with functional optical coherence tomography, *Scientific Reports*. 8 (2018) 1–12.
- [5] G. Ritter, L.S. Cohen, C. Williams, E.C. Richards, L.J. Old, S. Welt, Serological analysis of human anti-human antibody responses in colon cancer patients treated with repeated doses of humanized monoclonal antibody a33, *Cancer Research*. 61 (2001) 6851–6859.
- [6] E.M. Roth, A.C. Goldberg, A.L. Catapano, A. Torri, G.D. Yancopoulos, N. Stahl, A. Brunet, G. Lecorps, H.M. Colhoun, Antidrug antibodies in patients treated with alirocumab, (2017).
- [7] J.D. Jones, H.E. Ramser, A.E. Woessner, K.P. Quinn, In vivo multiphoton microscopy detects longitudinal metabolic changes associated with delayed skin wound healing, *Communications Biology*. 1 (2018) 1–8.
- [8] M.C. Skala, A.N. Fontanella, L. Lan, J.A. Izatt, M.W. Dewhirst, Longitudinal optical imaging of tumor metabolism and hemodynamics, *Journal of Biomedical Optics*. 15 (2010) 011112.
- [9] G.J. Greening, K.P. Miller, C.R. Spainhour, M.D. Cato, T.J. Muldoon, Effects of isoflurane anesthesia on physiological parameters in murine subcutaneous tumor allografts measured via diffuse reflectance spectroscopy, *Biomedical Optics Express*. 9 (2018) 2871–2886.
- [10] T.T. Sio, P.J. Atherton, B.J. Birkhead, D.J. Schwartz, J.A. Sloan, D.K. Seisler, J.A. Martenson, C.L. Loprinzi, P.C. Griffin, R.F. Morton, others, Repeated measures analyses of dermatitis symptom evolution in breast cancer patients receiving radiotherapy in a phase 3 randomized trial of mometasone furoate vs placebo (N06C4 [alliance]), *Supportive Care in Cancer*. 24 (2016) 3847–3855.
- [11] J. Kamstra, P. Dijkstra, M. Van Leeuwen, J. Roodenburg, J. Langendijk, Mouth opening in patients irradiated for head and neck cancer: A prospective repeated measures study, *Oral Oncology*. 51 (2015) 548–555.
- [12] E.-J. Wagenmakers, M. Lee, T. Lodewyckx, G.J. Iverson, Bayesian versus frequentist inference, in: *Bayesian Evaluation of Informative Hypotheses*, Springer, 2008: pp. 181–207.

- [13] R. Gueorguieva, J.H. Krystal, Move over anova: Progress in analyzing repeated-measures data and its reflection in papers published in the archives of general psychiatry, *Archives of General Psychiatry*. 61 (2004) 310–317.
- [14] P. Schober, T.R. Vetter, Repeated measures designs and analysis of longitudinal data: If at first you do not succeed—try, try again, *Anesthesia and Analgesia*. 127 (2018) 569.
- [15] J. Pinheiro, D. Bates, *Mixed-effects models in s and s-plus*, Springer Science & Business Media, 2006.
- [16] K. Vishwanath, H. Yuan, W.T. Barry, M.W. Dewhirst, N. Ramanujam, Using optical spectroscopy to longitudinally monitor physiological changes within solid tumors, *Neoplasia*. 11 (2009) 889–900.
- [17] B. Dennis, J.M. Ponciano, M.L. Taper, S.R. Lele, Errors in statistical inference under model misspecification: Evidence, hypothesis testing, and aic, *Frontiers in Ecology and Evolution*. 7 (2019) 372.
- [18] B. Wang, Z. Zhou, H. Wang, X.M. Tu, C. Feng, The p-value and model specification in statistics, *General Psychiatry*. 32 (2019).
- [19] C. Liu, T.P. Cripe, M.-O. Kim, Statistical issues in longitudinal data analysis for treatment efficacy studies in the biomedical sciences, *Molecular Therapy*. 18 (2010) 1724–1730.
- [20] L.G. Halsey, D. Curran-Everett, S.L. Vowler, G.B. Drummond, The fickle p value generates irreproducible results, *Nature Methods*. 12 (2015) 179–185.
- [21] H. Abdi, Holm’s sequential bonferroni procedure, *Encyclopedia of Research Design*. 1 (2010) 1–8.
- [22] S. Nakagawa, A farewell to bonferroni: The problems of low statistical power and publication bias, *Behavioral Ecology*. 15 (2004) 1044–1045.
- [23] A. Gelman, J. Hill, M. Yajima, Why we (usually) don’t have to worry about multiple comparisons, *Journal of Research on Educational Effectiveness*. 5 (2012) 189–211.
- [24] C. Albers, The problem with unadjusted multiple and sequential statistical testing, *Nature Communications*. 10 (2019) 1–4.
- [25] C. Ugrinowitsch, G.W. Fellingham, M.D. Ricard, Limitations of ordinary least squares models in analyzing repeated measures data, *Medicine and Science in Sports and Exercise*. 36 (2004) 2144–2148.
- [26] H. Huynh, L.S. Feldt, Estimation of the box correction for degrees of freedom from sample data in randomized block and split-plot designs, *Journal of Educational Statistics*. 1 (1976) 69–82.
- [27] S.W. Greenhouse, S. Geisser, On methods in the analysis of profile data, *Psychometrika*. 24 (1959) 95–112.
- [28] N. Haverkamp, A. Beauducel, 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 (2017) 1841.
- [29] H. Keselman, J. Algina, R.K. Kowalchuk, The analysis of repeated measures designs: A review, *British Journal of Mathematical and Statistical Psychology*. 54 (2001) 1–20.
- [30] J. Charan, N. Kantharia, How to calculate sample size in animal studies?, *Journal of Pharmacology & Pharmacotherapeutics*. 4 (2013) 303.
- [31] D.J. Barr, R. Levy, C. Scheepers, H.J. Tily, Random effects structure for confirmatory hypothesis testing: Keep it maximal, *Journal of Memory and Language*. 68 (2013) 255–278.
- [32] N.L. Rose, H. Yang, S.D. Turner, G.L. Simpson, An assessment of the mechanisms for the transfer of lead and mercury from atmospherically contaminated organic soils to lake sediments with particular reference to Scotland, UK, *Geochimica et Cosmochimica Acta*. 82 (2012) 113–135.
- [33] E.J. Pedersen, D.L. Miller, G.L. Simpson, N. Ross, Hierarchical generalized additive models in ecology: An introduction with mgcv, *PeerJ*. 7 (2019) e6876.
- [34] G.L. Simpson, Modelling palaeoecological time series using generalised additive models, *Frontiers in Ecology and Evolution*. 6 (2018) 149.

- [35] L. Yang, G. Qin, N. Zhao, C. Wang, G. Song, Using a generalized additive model with autoregressive terms to study the effects of daily temperature on mortality, *BMC Medical Research Methodology*. 12 (2012) 165.
- [36] N. Beck, S. Jackman, Beyond linearity by default: Generalized additive models, *American Journal of Political Science*. (1998) 596–627.
- [37] S.N. Wood, *Generalized additive models: An introduction with r*, CRC press, 2017.
- [38] R Core Team, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria, 2020. <https://www.R-project.org/>.
- [39] S.N. Wood, N., Pya, B. Säfken, Smoothing parameter and model selection for general smooth models (with discussion), *Journal of the American Statistical Association*. 111 (2016) 1548–1575.
- [40] B.T. West, K.B. Welch, A.T. Galecki, *Linear mixed models: A practical guide using statistical software*, CRC Press, 2014.
- [41] R.D. Wolfinger, Heterogeneous variance: Covariance structures for repeated measures, *Journal of Agricultural, Biological, and Environmental Statistics*. (1996) 205–230.
- [42] R.E. Weiss, *Modeling longitudinal data*, Springer Science & Business Media, 2005.
- [43] S. Geisser, S.W. Greenhouse, others, An extension of box’s results on the use of the F distribution in multivariate analysis, *The Annals of Mathematical Statistics*. 29 (1958) 885–891.
- [44] S.E. Maxwell, H.D. Delaney, K. Kelley, *Designing experiments and analyzing data: A model comparison perspective*, Routledge, 2017.
- [45] G. Molenberghs, H. Thijs, I. Jansen, C. Beunckens, M.G. Kenward, C. Mallinckrodt, R.J. Carroll, Analyzing incomplete longitudinal clinical trial data, *Biostatistics*. 5 (2004) 445–464.
- [46] J. Scheffer, *Dealing with missing data*, (2002).
- [47] R.F. Potthoff, G.E. Tudor, K.S. Pieper, V. Hasselblad, Can one assess whether missing data are missing at random in medical studies?, *Statistical Methods in Medical Research*. 15 (2006) 213–234.
- [48] Y. Ma, M. Mazumdar, S.G. Memtsoudis, Beyond repeated-measures analysis of variance: Advanced statistical methods for the analysis of longitudinal data in anesthesia research, *Regional Anesthesia & Pain Medicine*. 37 (2012) 99–105.
- [49] J. Pinheiro, D. Bates, S. DebRoy, D. Sarkar, R Core Team, *nlme: Linear and nonlinear mixed effects models*, 2020. <https://CRAN.R-project.org/package=nlme>.
- [50] J.A. Nelder, R.W. Wedderburn, Generalized linear models, *Journal of the Royal Statistical Society: Series A (General)*. 135 (1972) 370–384.
- [51] T. Hastie, R. Tibshirani, Generalized additive models: Some applications, *Journal of the American Statistical Association*. 82 (1987) 371–386.
- [52] T.J. Hefley, K.M. Broms, B.M. Brost, F.E. Buderman, S.L. Kay, H.R. Scharf, J.R. Tipton, P.J. Williams, M.B. Hooten, The basis function approach for modeling autocorrelation in ecological data, *Ecology*. 98 (2017) 632–646.
- [53] E.J. Wegman, I.W. Wright, Splines in statistics, *Journal of the American Statistical Association*. 78 (1983) 351–365.

A Simulation

A.1 Compound symmetry and independent errors in linear and quadratic responses

This section simulated linear and quadratic data in the same manner as in Section 2.5. The linear simulations using Figure 3 show in panels A and D the simulated mean responses and individual datapoints. Panels C and G show a visual interpretation of “correlation” in the responses: In panel C, subjects that have a value of the random error ε either above or below the mean group response are more likely to have other observations that follow the same trajectory, thereby demonstrating correlation in the response. In panel G, because the errors are independent, there is no expectation that responses are likely to follow a similar pattern. Panels D and H show the predictions from the rm-ANOVA model.

```
## Example with linear response
example <- function(n_time = 6,
                    fun_type = "linear",
                    error_type = "correlated") {

  if (!(fun_type %in% c("linear", "quadratic")))
    stop("'fun_type' must be either 'linear', or 'quadratic'")
  if (!(error_type %in% c("correlated", "independent")))
    stop("'error_type' must be either 'correlated', or 'independent'")

  library(tidyverse)
  library(mvnfast)
  library(nlme)
  library(ggsci)

  x <- seq(1,6, length.out = n_time)
  mu <- matrix(0, length(x), 2)
  # linear response
  if (fun_type == "linear") {
    mu[, 1] <- - (0.25*x)+2
    mu[, 2] <- 0.25*x+2
  } else {
    # nonlinear response

    mu[, 1] <- -(0.25 * x^2) +1.5*x-1.25
    mu[, 2] <- (0.25 * x^2) -1.5*x+1.25
  }
  # matplot(mu, type = 'l')

  y <- array(0, dim = c(length(x), 2, 10))
  errors <- array(0, dim = c(length(x), 2, 10))

  if (error_type == "independent") {
    ## independent errors
    for (i in 1:2) {
      for (j in 1:10) {
        errors[, i, j] <- rnorm(6, 0, 0.25)
        y[, i, j] <- mu[, i] + errors[, i, j]
      }
    }
  }
}
```

```

} else {
  for (i in 1:2) {      # number of treatments
    for (j in 1:10) {   # number of subjects
      # compound symmetry errors
      errors[, i, j] <- rmvn(1, rep(0, length(x)), 0.1 * diag(6) + 0.25 * matrix(1, 6, 6))
      y[, i, j] <- mu[, i] + errors[, i, j]
    }
  }
}

## subject random effects

## visualizing the difference between independent errors and compound symmetry
## why do we need to account for this -- overly confident inference

dimnames(y) <- list(time = x, treatment = 1:2, subject = 1:10)
dimnames(errors) <- list(time = x, treatment = 1:2, subject = 1:10)
dimnames(mu) <- list(time = x, treatment = 1:2)
dat <- as.data.frame.table(y, responseName = "y")
dat_errors <- as.data.frame.table(errors, responseName = "errors")
dat_mu <- as.data.frame.table(mu, responseName = "mu")
dat <- left_join(dat, dat_errors, by = c("time", "treatment", "subject"))
dat <- left_join(dat, dat_mu, by = c("time", "treatment"))
dat$time <- as.numeric(as.character(dat$time))
dat <- dat %>%
  mutate(subject = factor(paste(subject, treatment, sep = "-")))

## repeated measures ANOVA in R
fit_lm <- lm(y ~ time + treatment + time * treatment, data = dat)
dat$preds_lm <- predict(fit_lm)

fit_lme <- lme(y ~ treatment + time + treatment:time,
              data = dat,
              random = ~ 1 | subject,
              correlation = corCompSymm(form = ~ 1 | subject)
)

pred_dat <- expand.grid(
  treatment = factor(1:2),
  time = unique(dat$time)
)

dat$y_pred <- predict(fit_lme)

return(list(
  dat = dat,
  pred_dat = pred_dat,
  fit_lm = fit_lm,

```

```

    fit_lme = fit_lme
  })
}

plot_example <- function(sim_dat) {
  library(patchwork)
  ## Plot the simulated data
  p1 <- sim_dat$dat %>%
    ggplot(aes(x = time, y = y, group = treatment, color = treatment)) +
    geom_point(show.legend=FALSE) +labs(y='response')+
    geom_line(aes(x = time, y = mu, color = treatment),show.legend=FALSE) +
    theme_classic() +
    #ggtitle("Simulated data with true response function")+
    theme(plot.title = element_text(size = 30,
                                     face = "bold"),
          text=element_text(size=30))+
    scale_color_aaes()

  p2 <- sim_dat$dat %>%
    ggplot(aes(x = time, y = y, group = subject, color = treatment)) +
    geom_line(aes(size = "Subjects"),show.legend = FALSE) +
    # facet_wrap(~ treatment) +
    geom_line(aes(x = time, y = mu, color = treatment, size = "Simulated Truth"), lty = 1,show.legend =
    scale_size_manual(name = "Type", values=c("Subjects" = 0.5, "Simulated Truth" = 3)) +
    #ggtitle("Simulated data\nIndividual responses with population mean") +
    theme_classic()+
    theme(plot.title = element_text(size = 30,
                                     face = "bold"),
          text=element_text(size=30))+
    scale_color_aaes()

  p3 <- sim_dat$dat %>%
    ggplot(aes(x = time, y = errors, group = subject, color = treatment)) +
    geom_line(show.legend=FALSE) +labs(y='errors')+
    theme_classic()+
    # facet_wrap(~ treatment) +
    #ggtitle("Simulated errors") +
    theme(plot.title = element_text(size = 30,
                                     face = "bold"),
          text=element_text(size=30))+
    scale_color_aaes()

  p4 <- ggplot(sim_dat$dat, aes(x = time, y = y, color = treatment)) +
    geom_point()+labs(y='response')+
    geom_line(aes(y = predict(sim_dat$fit_lme), group = subject, size = "Subjects")) +
    geom_line(data = sim_dat$pred_dat, aes(y = predict(sim_dat$fit_lme, level = 0, newdata = sim_dat$pr
    scale_size_manual(name = "Predictions", values=c("Subjects" = 0.5, "Population" = 3)) +
    theme_classic() +
    #ggtitle("Fitted Model")+
    theme(plot.title = element_text(size = 30,
                                     face = "bold"),
          text=element_text(size=30))+

```



```

    scale_color_aaas()

    return((p1+p3+p2+p4)+plot_layout(nrow=1)+plot_annotation(tag_levels = 'A'))
  }

txt<-18
A1<-plot_example(example(fun_type = "linear", error_type = "correlated"))

B1<-plot_example(example(fun_type = "linear", error_type = "independent"))

C1<-plot_example(example(fun_type = "quadratic", error_type = "correlated"))

D1<-plot_example(example(fun_type = "quadratic", error_type = "independent"))

```

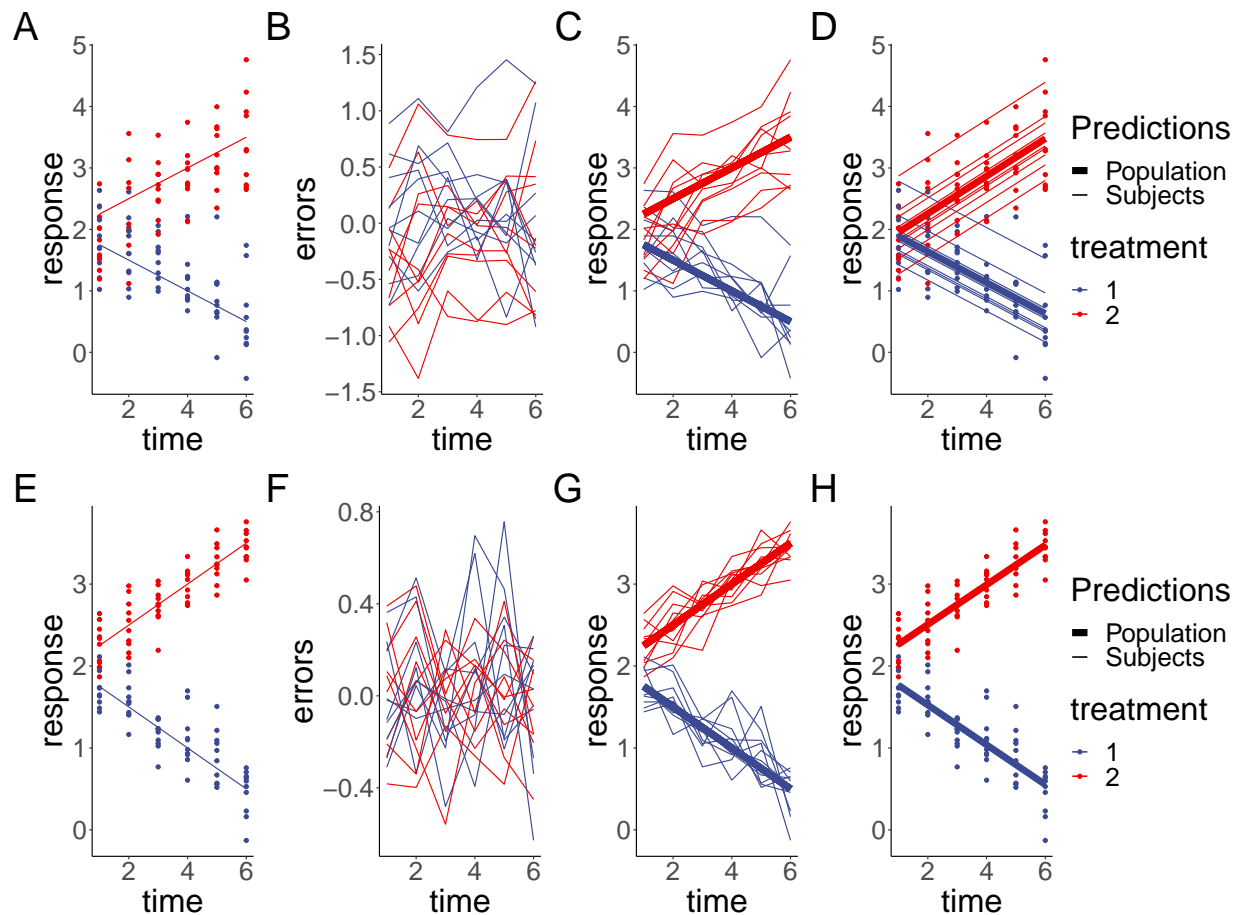


Figure 3: **Simulated linear responses from two groups with correlated (top row) or independent (bottom row) errors using a rm-ANOVA model. A, C: Simulated data with known mean response (linear or quadratic, thin lines) and individual responses (points) showing the dispersion of the data. B, D: Estimations from the rm-ANOVA model for the mean group response (linear or quadratic). Thick lines are the predicted mean response per group, thin lines are the random effects for each subject and points represent the original raw data. The rm-ANOVA model does not pick the trend of the quadratic data**

For the quadratic response case, Figure 4 shows the simulated responses using compound symmetry and independent errors.

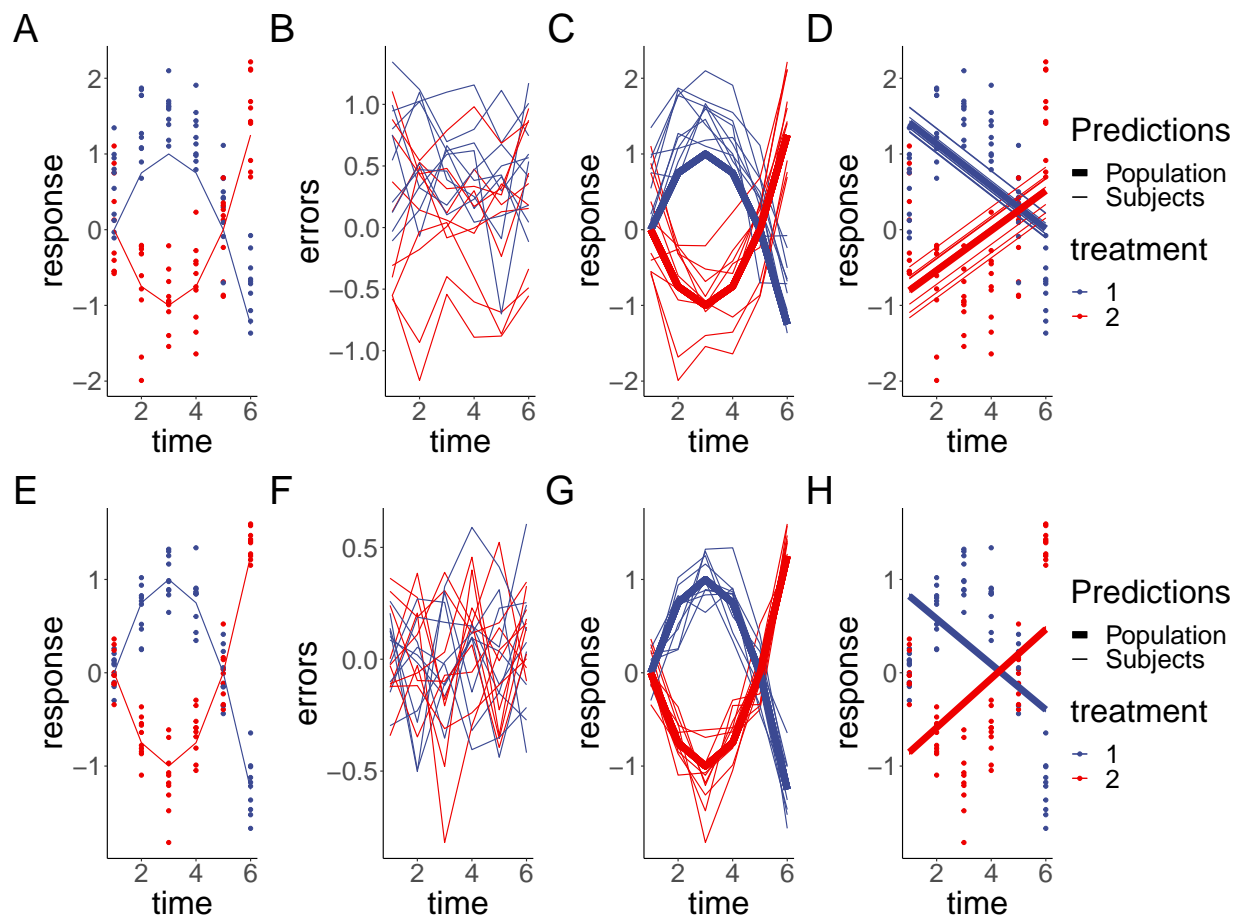


Figure 4: **Simulated quadratic responses from two groups with a rm-ANOVA model fitted. A,E: Simulated data with known mean response (lines) and individual responses (points) showing the dispersion of the data. B,F: Generated errors showing the difference in the behavior of correlated and independent errors. C,G: Simulated known response per group (thick lines) with individual trajectories (thin lines), note that subjects with observations in the area above the mean response tend to stay in that region through the timeline. D,H: Estimations from the rm-ANOVA model for the mean group response. Thick lines are the predicted mean response per group, thin lines are the random effects for each subject and points represent the original raw data.**