

# Introduction

In standard undergraduate curricula, there is a strong focus on cross sectional data, and thus no emphasis on how time-sequence data is analyzed. However, a significant portion of data that we encounter in the real world is dependent on time. If we want to track trends and changes over time, such as an effect of a certain drug on the body or growth of a company, longitudinal data and analysis will help us examine those points of interest. For example, the Chinese Longitudinal Healthy Longevity Survey from Duke University assessed physical and mental well-being of Chinese elders for over almost 2 decades and re-interviewed survivors every few year. This follow up in data collection allowed researchers to investigate the aging process over time and identify risk factors and causes leading up to death.

Not only can we observe change over time in individuals, but we can look at higher-level grouping, such as change in schools, counties, and organizations. It should be emphasized that only longitudinal data can capture changes within a subject or group; cross-sectional data contain responses that are captured at only one occasion that are then compared to other subjects. Ultimately, it cannot provide information about changes over time.

One key aspect of longitudinal data is that there needs to be repeated measurements of the same individuals across multiple periods of time. If there aren't repeated observations, then it is not possible to make any comparisons between two or more time points. Having repeated measurements of the same individual allows for removal of potential confounding effects, such as gender or socioeconomic status, from the analysis. Since we assume that these confounding variables are fixed effects that do not vary from measurement to measurement, all changes from an individual cannot be attributed to these effects.

The measure that captures the observed changes within an individual is referred to as a response trajectory. There are different ways of comparing response trajectories. For example, it is possible to compare the post-treatment vs baseline changes across multiple treatment groups, or it is also possible to compare the rate of change. The method chosen depends on the specific question of the study.

Apart from comparing just the response trajectories, it is also of interest to compare individual differences in the relationship between covariates and the response trajectory. This can be captured using various different statistical models. The choice of model depends on several characteristics of the data.

## Characteristics of longitudinal data

While the only requirement of longitudinal data is that there is more than one observation for a given individual, there are other components that affect the model chosen. Data can be unbalanced or balanced: *balanced* data refers to when all individuals have the same number of repeated measurements taken at the same occasions. In addition, data can also be missing, resulting in automatically unbalanced data. This affects the accuracy of how changes over time are analyzed depending on if there are any patterns to the missing data or not.

Another unique characteristic of longitudinal data is that repeated measurements of each individual are typically positively correlated. This feature violates conditions of other common statistical methods such as linear regression, where measurements are assumed to be independent. This positive correlation allows for more accurate estimates of the model coefficients and response trajectories since there is reduced uncertainty knowing that a previous measurement can help predict the next one.

In longitudinal analysis, a covariance matrix is calculated for each individual and all of their measurements. The diagonals of this matrix represent the variance of each of the measurements, which are not constant over time. The off-diagonals of the matrix are non-zero to account for the lack of independence between measurements, and are usually not constant because correlations between measurements decrease over time. While these values are rarely 0, they are also rarely 1. There are different covariance pattern structures that are imposed that account for these features.

These features of the covariance of longitudinal data serve as the underlying premise to the idea that variation can be separated into three distinct parts: 1) between-individual variation, 2) within-individual variation, and 3) measurement error.

Between-individual variation helps explain why measurements from the same individual are more likely to be positively correlated than measurements to a different individual. Within-individual variation helps explain why correlations decrease with increasing time differences, and measurement error explains why correlations are never one. These three types of variation may contribute to total variation in unequal amounts, but may not need to be differentiated depending on the type of longitudinal analysis desired.

## Notation

Throughout the rest of the text, we will use a standard set of notation for all parameters and variables.  $Y_{ij}$  represents the response variable for the  $i^{th}$  individual at the  $j^{th}$  measurement. When we have repeated  $n_i$  measurements for an individual, we can construct a vector,

$$Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{pmatrix}.$$

We use  $\mu_{ij}$  as the conditional mean response at the  $j^{th}$  measurement, where conditional entails a dependence of the mean response on the covariates.

## Estimation and Inference

Regression coefficient values  $\beta$  and the covariance matrix  $\Sigma_i$  can be estimated using maximum likelihood estimation, which identifies values of  $\beta$  and  $\Sigma_i$  that maximize the joint probability of the response variable occurring based on the observed data; the probability is known as the likelihood function. These values are estimates that are denoted by  $\hat{\beta}$  and  $\hat{\Sigma}_i$ . When observations are independent of one another, maximizing the likelihood function for  $\beta$  is equivalent to finding a value of  $\hat{\beta}$  that minimizes the sum of the squares of the residuals. However, since there are repeated measurements of each individual that are not independent of one another we use the generalized least squares (GLS) estimator:

$$\hat{\beta} = \{\Sigma_{i=1}^N (X_i' \Sigma_i^{-1} X_i)\}^{-1} \Sigma_{i=1}^N (X_i' \Sigma_i^{-1} y_i).$$

In addition, the sampling distribution of  $\hat{\beta}$  has mean  $\beta$  and covariance:

$$\hat{Cov}(\hat{\beta}) = \{\Sigma_{i=1}^N (X_i' \Sigma_i^{-1} X_i)\}^{-1}.$$

The GLS estimator assumes that  $\Sigma_i$  is known. However, since this isn't usually the case, we can substitute  $\Sigma_i$  with a maximum likelihood estimate  $\hat{\Sigma}_i$ . It can be shown that the properties of  $\hat{\beta}$  still hold using an estimate of the covariance.

While the maximum likelihood estimate of  $\Sigma_i$  is adequate, a modified method known as restricted maximum likelihood (REML) estimation is suggested to reduce bias in finite samples. The bias originates from the fact that  $\beta$  itself is also estimated from data, but is not accounted for when estimating covariance. In REML estimation of  $\Sigma_i$ ,  $\beta$  is removed from the likelihood function. This REML estimation of  $\Sigma_i$  can be used in the GLS estimator for  $\hat{\beta}$  mentioned above, and is recommended in place of the ML estimator.

Now that we have estimates for  $\beta$ , we can make inferences through construction of confidence intervals and hypothesis testing. For example, using the ML estimate  $\hat{\beta}$  and  $\hat{Cov}(\hat{\beta})$ , we can construct a Wald statistic to test for significance of  $\hat{\beta}_k$ :

$$Z = \frac{\hat{\beta}_k}{\sqrt{\hat{Var}(\hat{\beta}_k)}}.$$

One crucial assumption when conducting inference using the ML estimate for  $\beta$

## Linear models for longitudinal data

As mentioned previously, there are multiple ways to model longitudinal data. When the response variable is continuous, we can consider a model that relates the mean response and the covariates in a linear way. In a linear model all components can be represented using vectors and matrices. The most general form of the linear model can be represented as:

$$E(Y|X_i) = X_i\beta$$

, where  $\beta$  is a vector of regression coefficients and  $X_i$  is a vector of covariates. We will discuss three methods for linear models: 1) response profile analysis, 2) parametric time model, 3) linear mixed effect model.

### Response profile analysis

In response profile analysis, we allow for arbitrary patterns in the mean response over time. A sequence of means over time is known as the mean response profile. The main goal of this analysis is to identify differences in pattern of change in mean response profile among 2 or more groups. This method requires that the data be balanced.

There are three effects of interest when analyzing response profiles in longitudinal analysis: 1. *group*  $\times$  *time* interaction effect (are the mean response profiles different in groups over time?) 2. time effect (assuming mean response profiles are parallel between groups, are the means changing over time?) 3. Group effect (do the mean response profiles differ?)

However, the first question is the primary interest. The goal is to find whether the change in mean response over time differs across groups.

To test for significance of the *group*  $\times$  *time* effect, we have a null hypothesis that the difference in means between the  $n$  groups is constant over time, which in other words entails that mean response profiles between the groups have parallel slopes. We can implement the general linear model  $\mu_i = X_i\beta$  to test our hypotheses, using comparison of  $\beta$  slope parameters to determine whether there is a *group*  $\times$  *time* effect.

For example, to express the model for response profile analysis for  $G$  groups and  $n$  occasions of measurement, we have  $G \times n$  parameters for the  $G$  mean response profiles. For two groups measured at three occasions, we have 6 slope parameters. if  $\beta_1 - \beta_3$  represent slope parameters for mean responses in group 1 and  $\beta_4 - \beta_6$  represent slope parameters for mean responses in group 2, our null hypotheses would be that  $(\beta_1 - \beta_4) = (\beta_2 - \beta_5) = (\beta_3 - \beta_6)$ .

An unstructured covariance model is typically assumed for response profile analysis. “Unstructured” means that there is no explicit structure or pattern imposed on the covariance for the repeated measures, so each of the variances and covariance pairs are estimated using restricted maximum likelihood estimation (REML). For  $n$  repeated measures, there are  $n$  variances and  $n \times (n - 1)/2$  covariances to be estimated. In a study where there are 10 repeated measurements, there 55 total covariance parameters to be estimated, which can become computationally intensive.

One other aspect to consider when conducting analysis on mean response profiles is how to adjust for the baseline measurement. The baseline value is important when we want to calculate measures that compare mean response to the baseline. How we adjust depends on whether the study is randomized or observational. When the study is randomized and baseline measurement is taken before treatment assignment, the mean response at occasion 1 is independent of the group, and assumed to be equal. One possible method is to treat the baseline measurement as a covariate, and use response measurements 2 through  $n$  as the dependent measures. This is referred to as the analysis of covariance approach. Additionally, this method only works for randomized studies because using the baseline measurement as a covariate for observational studies may produce confounding effects. For an observational study, it is recommended to subtract the baseline response

to create a change score. For both types of longitudinal studies there are various methods to account for the baseline value, and should be considered carefully before implementing the method.

Overall, response profile analysis is a straightforward method in investigating differences between groups for longitudinal data. Since both the covariance and mean responses have no imposed structure, the analysis is more robust and immune to inaccurate results due to model misspecification. However, there are drawbacks as well. Response profile analysis does not consider time-order of the measurements and does not distinguish between between-individual variation and within-individual variation. In addition, it can only provide a broad analysis of whether there are differences across groups and time, but does not provide the amount of detail usually needed to answer research questions, such as how exactly measurements taken towards the end of the study compare to measurements taken at the beginning. In this method, time is treated as a categorical covariate rather than a continuous one. Another method that addresses the issue of examining time order of the data is parametric time models.

## Parametric Time Models

Parametric time models are able to capture time order of the data by fitting linear or quadratic curves to capture an increasing or decreasing pattern over time. Time is treated as a continuous covariate rather than a categorical one. In addition, unlike response profile analysis, parametric time models are able to handle unbalanced and missing data. Rather than fitting a complex and perfect model onto the observed mean response profile, parametric time models fit simple curves that produce covariate effects of greater power. This is because in mean response profile we are testing a wider range of hypotheses since we are looking for inequality between two groups; however, in parametric time models, we are testing more specifically whether the data follow a linear trend, which results in more power.

Additionally, while in the mean response profile analysis an unstructured covariance pattern is assumed, here there is flexibility in choice of the covariance model; there are several options such as Toeplitz or compound symmetric that impose various structures on the model. For example, a Toeplitz model:

$$Cov(Y_i) = \begin{pmatrix} 1 & \rho_1 & \rho_2 & \dots & \rho_{n-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_{n-2} \\ \rho_2 & \rho_1 & 1 & \dots & \dots \\ \rho_3 & & & & \end{pmatrix}$$

structures the covariance matrix such that any pair of responses that are equally separated in time have the same correlation.

It is possible to choose an unstructured covariance model as well, but can be computationally intense if there are a large number of measurements.

We can use parametric time models in two ways: through polynomial trends and linear spines.

## Polynomial Trends

Using polynomial trends such as linear or quadratic, we can model longitudinal data as a function of time. Linear trends are the most common and interpretable ways to model change in mean over time. In an example comparing a treatment group to a control group, we can fit a linear trend using the following equation:

$$E(Y_{ij}) = \beta_1 + \beta_2 Time_{ij} + \beta_3 Group_i + \beta_4 Time_{ij} \times Group_i.$$

If  $\beta_4 = 0$ , then the two groups do not differ in terms of changes in the mean response over time.

For quadratic trends, the changes in mean are no longer constant since the rate of change depends on the time. Thus, we fit an additional parameter to express the rate of change. Using the previous example of treatment vs. control group, we have the model:

$$E(Y_{ij}) = \beta_1 + \beta_2 Time_{ij} + \beta_3 Time_{ij}^2 + \beta_4 Group_i + \beta_5 Time_{ij} \times Group_i + \beta_6 Time_{ij}^2 \times Group_i.$$

As we can see from the models above, the inclusion of an additional parameter  $Time_{ij}^2$  changes the mean response rate. One problem that may arise from using quadratic trends is that there is collinearity between  $Time_{ij}$  and  $Time_{ij}^2$ , which can affect the estimation of  $\beta$ . To account for this, we can center the  $Time_{ij}$  variable around the mean time value for all individuals, instead of centering it around zero as done in normal analysis. For example if we have a set of times  $Time = 0, 1, 2, \dots, 10$ , then the mean time value is five. Thus time zero would be recentered as -5. The interpretation of the intercept changes to represent the mean response at that recentered mean time value.

## Linear splines

In instances where responses cannot be adequately fit by polynomial trends, such as when the responses fluctuate between increasing and decrease at different extents, we can employ a linear spline model. This model consists of piece-wise line segments that have unique slopes for a given set of time measurements. The point at which different line segments meet are called knots, and the number of knots depends on the context of the data and researcher discretion.

Drawing again from our treatment vs control group design, a linear model for the mean responses of the control group is:

$$E(Y_{ij}) = \beta_1 + \beta_2 Time_{ij} + \beta_3 (Time_{ij} - t^*)_+.$$

The  $()_+$  indicates a truncated line function and is positive when  $Time_{ij} - t^*$  is greater than 0, and otherwise is equal to 0. In this case, the function depends on the specified time  $t^*$ . If the mean response is before  $t^*$ , then the mean response is modeled by:

$$E(Y_{ij}) = \beta_1 + \beta_2 Time_{ij}.$$

If the mean response is after  $t^*$ , it is modeled by

$$E(Y_{ij}) = (\beta_1 - \beta_3 t^*) + (\beta_2 + \beta_3) Time_{ij}.$$

There are benefits to parametric models that make them a more appealing choice compared to response profile analysis. Parametric time models are able to capture time order, and can be used with unbalanced data. However, they do not differentiate between subject and within subject variation. If further analysis of individual variation is desired, linear mixed effects models can be employed.

## Linear Mixed Effects

In both response profile analysis and parametric time models, the regression parameters are considered to be universal for each population group. However, in instances where we want to account for heterogeneity within a population, we can use a linear mixed effects model and consider a subset of the regression parameters to be random. This model distinguishes between fixed effects, which are population characteristics shared by all individuals, and subject specific effects, also known as random effects, which pertain to each individual. These subject specific effects mean that parameters are random, which induces a structure onto the covariance model.

In addition, distinguishing between fixed and random effects allows for differentiation between within-subject and between-subject variation.

One example of the linear mixed effects model is the random intercept model, which is the simplest version of the linear mixed effects model:

$$Y_{ij} = X'_{ij}\beta + b_i + \epsilon_{ij}$$

This model is very similar to the general linear model with a few additions.  $b_i$  is the random subject effect and  $\epsilon$  is the measurement error. Both effects are random, with mean 0 and  $\text{Var}(b_i) = \sigma_b^2$ ,  $\text{Var}(\epsilon_{ij}) = \sigma^2$ .

$X'_{ij}\beta$  is the population mean, and  $b_i$  represents the differing subject effect that is unique to each individual.  $b_i$  is interpreted as how the subject deviates from the population mean while accounting for covariates.

As mentioned previously, the random effects are responsible for inducing a structure on the covariance model. This structure is not to be confused with the covariance structures that can be chosen when using parametric time models. For a given individual, it can be shown that variance of each response is:

$$\text{Var}(Y_{ij}) = \sigma_b^2 + \sigma^2$$

and the covariance between two measurements  $Y_{ij}$  and  $Y_{ik}$  is equal to  $\sigma_b^2$ . The resulting covariance matrix

$$\begin{pmatrix} \sigma_b^2 + \sigma^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma^2 & \sigma_b^2 + \sigma^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma^2 & \dots & \sigma_b^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 + \sigma^2 \end{pmatrix}$$

implies correlation between measurements, and also highlights the role played by the random effects in determining the covariance.

Extending beyond the random intercept model, multiple random effects can be incorporated.

A linear mixed effects model can be expressed as

$$Y_i = X_i\beta + Z_ib_i + \epsilon_i.$$

Where:  $\beta$  is a  $p \times 1$  vector of fixed effects  $b_i$  is a  $q \times 1$  vector of random effects  $X_i$  is a  $n \times p$  matrix of covariates  $Z_i$  is a  $n \times q$  matrix of covariates

The subset of regression covariates that vary randomly are found in  $Z_i$ . We assume that  $b_i$  comes from a multivariate normal distribution with mean 0 and covariance matrix  $G$ . We also assume that  $\epsilon_i$  are independent of  $b_i$ , and come from multivariate normal distribution with mean 0 and covariance matrix  $R_i$ .

The covariance of  $Y_i$  can be modeled by

$$\text{Cov}(Z_ib_i) + \text{Cov}(\epsilon_i) = Z_iGZ_i' + R_i.$$

This model, which outlines a distinction between  $G$  and  $R_i$ , allows for separate analysis of between subject and within subject variation. Unlike other covariance models, in linear mixed effects models the covariance is a function of the times of measurement. This allows for unbalanced data to be used for the model since each individual can have their unique set of measurement times. Lastly, the model allows for variance and covariance to change as a function of time. To illustrate, consider the following model:

In an example where individuals can vary both in their baseline response and their rate of change, we have:

$$Y_i = X_i\beta + Z_ib_i + \epsilon_i,$$

where both  $X_i$  and  $Z_i = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \dots & \dots \\ 1 & t_{in} \end{pmatrix}$ . For the  $i^{th}$  subject at the  $j^{th}$  measurement, the equation is as follows:

$$Y_{ij} = \beta_1 + \beta_2 t_{ij} + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij}.$$

If  $\text{Var}(b_{1i}) = g_{11}$ ,  $\text{Var}(b_{2i}) = g_{22}$ , and  $\text{Cov}(b_{1i}, b_{2i}) = g_{12}$  where these three components represent the  $G$  covariance for  $b_i$ , then it can be shown that  $\text{Cov}(Y_{ij}, Y_{ik}) = g_{11} + (t_{ij} + t_{ik})g_{12} + t_{ij}t_{ik}g_{22}$ .

Here in the covariance matrix we can see the dependence of the covariance on time. In this example there are four covariance parameters that arise from the two random effects of intercept and time. The number of covariance parameters is represented by  $q \times (q + 1)/2 + 1$ , where  $q$  is the number of random effects. To choose the most optimal model for covariance, we compare two nested models, one with  $q + 1$  random effects and one with  $q$  random effects. We use the likelihood ratio test to make a decision for which model to use.

One additional analysis that is possible with linear mixed effects models is predicting subject-specific responses. Given that  $b_i$  is a random variable, we can predict it using:

$$E(b_i|Y_i) = GZ_i(\Sigma)_i^{-1}(Y_i - X_i\hat{\beta}).$$

Because the covariance of  $Y_i$  is unknown, we can estimate both  $G$  and  $(\Sigma)_i^{-1}$  using REML, creating  $\hat{b}_i$ , also known as the empirical best linear unbiased prediction (BLUP). Thus, the equation for predicting the response profile is:

$$\hat{Y}_i = X_i\hat{\beta} + Z_i\hat{b}_i.$$

This equation to estimate the mean response profile can be extended to incorporate  $R_i$ , which represents within-subject variability. From this extension, we see that the equation and the empirical BLUP account for the weighting of both the within-subject variability and between-subject variability. If there is more within-subject variability, then more weight is assigned to  $X_i\hat{\beta}$ , the population mean response profile, in comparison to the subject's individual responses, and vice versa.

## Choosing the best model

After presenting three methods of evaluating longitudinal data, the natural question arises of how to choose the most appropriate model. While there is no definite correct answer, there are several factors to consider. If data are unbalanced, response profile analysis should not be considered; rather, parametric time model or linear mixed effect model would be more optimal. If time order is important to the analysis, then only parametric time model and linear mixed effect model should be used. If there is a need to distinguish between the two types of variation that can occur, then only linear mixed effect models are appropriate. The model should ultimately be chosen based on the characteristics and constraints of the data, as well as the specificity of the research question at hand.

## Conclusion

Longitudinal analysis is a valuable method to analyze changes over time. It is important to understand the unique characteristics that come with this analysis and to choose the best model that can capture the salient patterns that arise from the data. In subsequent chapters we will dive more deeply into how inference in longitudinal analysis is affected when sample sizes are not efficient through both simulation and application.