

# Table of Contents

<b>Abstract</b> . . . . .	i
<b>Acknowledgments</b> . . . . .	iii
<b>List of Tables</b> . . . . .	v
<b>List of Figures</b> . . . . .	vii
<b>Chapter 1: Introduction</b> . . . . .	1
1.1 Characteristics of Longitudinal Data . . . . .	2
1.2 Estimation and Inference . . . . .	4
1.3 Linear Models for Longitudinal Data . . . . .	5
1.3.1 Notation and Distribution Assumptions . . . . .	5
1.3.2 Response Profile Analysis . . . . .	7
1.3.3 Parametric Time Models . . . . .	8
1.3.4 Polynomial Trends . . . . .	10
1.3.5 Linear Mixed Effects . . . . .	11
1.4 Choosing the Best Model . . . . .	15
1.5 Conclusion . . . . .	15
<b>Chapter 2: Degrees of Freedom Methods and Simulation Study</b> . . .	17
2.1 Inference . . . . .	17
2.1.1 Inference in Small Sample Sizes . . . . .	18
2.2 Satterthwaite . . . . .	19

2.3	Kenward-Roger . . . . .	21
2.4	Other Methods . . . . .	22
2.5	Existing Literature . . . . .	22
2.6	Goals of This Study: . . . . .	24
2.7	Simulation Set Up: . . . . .	24
2.7.1	Generating Data: Sample Dize . . . . .	24
2.7.2	Generating Data: Fixed Effects . . . . .	25
2.7.3	Generating Data: Random effects . . . . .	25
2.7.4	Linear Mixed Effects Model . . . . .	27
2.8	Evaluating and Results . . . . .	28
2.8.1	Exponential Distribution . . . . .	31
2.9	Lognormal Distribution . . . . .	33
2.10	Normal Distribution . . . . .	35
2.10.1	KR vs Satterthwaite . . . . .	37
2.11	Discussion . . . . .	38
<b>Chapter 3: Application</b>	. . . . .	<b>43</b>
3.1	Application to Longitudinal Study about Children's Health . . . . .	43
3.1.1	Background . . . . .	43
3.1.2	Exploration . . . . .	44
3.1.3	Linear Mixed Model . . . . .	46
3.1.4	Intercept Only Model . . . . .	46
3.1.5	Intercept and Random Slope . . . . .	48
3.2	Discussion . . . . .	50
<b>Conclusion</b>	. . . . .	<b>53</b>
<b>Appendix A: The First Appendix</b>	. . . . .	<b>55</b>
A.1	All Code: . . . . .	55

# Chapter 1 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 two decades and re-interviewed survivors every few years (Zeng, Vaupel, Xiao, Liu, & Zhang, 2017). 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 may also be of interest to compare individual differences in the relationship between covariates and the response trajectory. This can be captured using various statistical models. The choice of model depends on several characteristics of the data.

## 1.1 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 characteristics 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. For example, a study measuring the weights of children over time may be unbalanced if some children are measured at 5 and 10 years old, while others are measured at 5, 7 and 15 years old, since measurements are not taken at the same time points. If children drop out of the study as they get older, this will result in missing data

→ Missing data are not a unique characteristic of longitudinal data, so this feels introduced incorrectly. It does make longitudinal data unbalanced, but imbalance can occur in other ways. Maybe re-frame to clarify?

points.

Another unique characteristic of longitudinal data is that repeated measurements within each individual are typically positively correlated, while measurements between individuals are independent (Fitzmaurice, Laird, & Ware, 2011). 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.

We can capture the associations among repeated measures within each individual by constructing a covariance matrix. In longitudinal analysis, a covariance matrix is calculated for each individual and all of their measurements. The diagonal element  $\sigma^2$  of this matrix represents the variance of each of the measurements, which may not be constant over time for longitudinal data. The off-diagonal elements of the matrix are non-zero to account for the lack of independence between measurements, and are usually not constant because correlations between measurements tend to decrease over time. While these values are rarely 0, they are also rarely 1 (Fitzmaurice et al., 2011). There are different covariance pattern structures that are imposed that account for these features.

These features of the covariance of longitudinal data can be explained when we separate the total variation 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 between different individuals. Within-individual variation helps explain why correlations between repeated measurements decrease with increasing time differences, and measurement

error along with individual variations 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.

Ch 2

## 1.2 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: cite FLW!

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

Should this be N?  
See Notation section  
Comments

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

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

$\overset{\text{no hat}}{Cov(\hat{\beta})} = \text{Cov}(\hat{\beta})$  with  $\hat{\Sigma}_i$  plugged in

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

$$\hat{Cov}(\hat{\beta}) = \sum_{i=1}^n (X_i' \hat{\Sigma}_i^{-1} X_i)$$

What properties?  
Other include properties of  $\hat{\beta}$   
before this paragraph  
or delete this line.

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  $\hat{\beta}$  of  $\Sigma_i$  can be used in the GLS estimator for  $\hat{\beta}$  mentioned above, and is recommended in place of the ML estimator.

*Split into separate inference section* 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  $\hat{\beta}_k$ :

*Xwidehat{Cov}*

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

*\widehat{\text{Var}}*

In later chapters we will explore how inference may be impacted in smaller sample sizes.

## 1.3 Linear Models for Longitudinal Data

### 1.3.1 Notation and Distribution Assumptions

Throughout the rest of the text, we will use a standard set of notation.  $Y_{ij}$  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,

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

*i = 1, ..., I?*  
*how are you representing # of individuals? i = 1, ..., I?*  
*total # of msmts?*  
 $N = \sum_{i=1}^I n_i$

$X_{ij} = \begin{bmatrix} x_{1,ij} \\ \vdots \\ x_{p-1,ij} \end{bmatrix}$   
 flip  
 flip  
 For each  $Y_{ij}$  response, there is a  $X_{ij}$  vector of covariates that has  $p$  rows representing whose first entry is 1.  
the number of covariates. Covariates in  $X_{ij}$  can be fixed or change over time.

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 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:

$$Y_i = X_i \beta + e_i,$$

in matrix form

$X_i = n_i \times p$   
 $\beta = p \times 1$

Check dimensions and notation

where  $X_i$  is a matrix representing the grouped collection of  $X_{ij}$  vectors, with each row representing a unique measurement, and  $p$  columns for each covariate that is associated with  $Y_i$ .  $\beta = \beta_1, \dots, \beta_p$  is a vector of regression coefficients that quantifies the relationship between the response and each covariate.  $e_i$  is a  $n_i \times 1$  vector of random errors for each measurement.

The linear model can be divided into a systematic component,  $X_i\beta$ , and a random component  $e_i$ . These two components contribute to the distribution assumptions of  $Y_i$ .  $Y_i$  is assumed to have a conditional multivariate normal distribution with mean response vector

$$E(Y_i|X_i) = X_i\beta,$$

which is the systematic component, and the covariance matrix

$$\Sigma_i = \text{Cov}(Y_i|X_i),$$

which captures the random variability of  $Y_i$ , the random component, and its role in shaping the overall distribution. In addition, this distribution is considered conditional because of its dependence on the covariates  $X_i$ .

What about  $\beta$ ?  
 Also, I would just write " $\beta$  is a  $p \times 1$  vector of regr coeffs..."

Can be more succinct with notation, e.g.

We assume

$$Y_i | X_i \sim MVN(\mu_i, \Sigma_i)$$

where  $\mu_i = E(Y_i|X_i) = X_i\beta$

$$\text{and } \Sigma_i = \text{Cov}(Y_i|X_i)$$

*Not adding much*

Oftentimes, the mean response vector  $E(Y_i|X_i) = X_i\beta$  is used for modeling longitudinal responses. In the following sections, we will discuss three specific methods of linear models: 1) response profile analysis, 2) parametric time model, and 3) linear mixed effects model.

### 1.3.2 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 profiles among two or more groups. This method requires that the data be balanced. *and treats time as a categorical variable.*

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?)

Analyzing response profiles can be modeled using

$$E(Y_i|X_i) = X_i\beta.$$

In an example where two groups with three measurements each are compared,  $\mu(1) = \beta_1, \beta_2, \beta_3$  represent regression coefficients for group 1, and  $\mu(2) = \beta_4, \beta_5, \beta_6$  represent regression coefficients for group 2. The group  $\times$  time interaction effect can be expressed as a null hypothesis in the form

$$H_0 : \beta_5 = \beta_6 = 0.$$

An unstructured covariance model is typically assumed for response profile anal-

ysis. “Unstructured” means that there is no explicit structure or pattern imposed on the covariance for the repeated measures. This is represented as

$$\Sigma := \text{Cov}(Y_i) = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix}.$$

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 are 55 total covariance parameters to be estimated, which can become computationally intensive.

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 needed to answer certain 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.

### 1.3.3 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. The same question, such as examining group  $\times$  time effect, requires  $n$  parameters in response profile analysis, but only requires one parameter in parametric time models.

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:

$$\Sigma_i = \sigma^2 \begin{pmatrix} 1 & \rho_1 & \rho_2 & \dots & \rho_{n-1} \\ \rho_1 & 1 & \rho_1 & \dots & \rho_{n-2} \\ \rho_2 & \rho_1 & 1 & \dots & \dots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho_{n-1} & \rho_{n-2} & \rho_{n-3} & \dots & 1 \end{pmatrix}$$

structures the covariance matrix such that any pair of responses that are equally separated in time have the same correlation. A compound symmetry model:

$$\Sigma_i \underset{\text{Cov}(Y_i)}{=} \sigma^2 \begin{pmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \rho & \rho & 1 & \dots & \rho \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \rho & \dots & 1 \end{pmatrix},$$

assumes constant variance and same correlation across all measurements of  $Y_i$ . It is

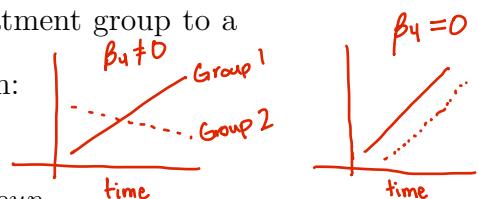
possible to choose an unstructured covariance model as well, but can be computationally intense if there are a large number of measurements.

We can use parametric time models in two ways: through polynomial trends and linear splines. More information about linear splines can be found in Fitzmaurice et al. (2011). *specify section*

#### 1.3.4 Polynomial Trends *remove subheading*

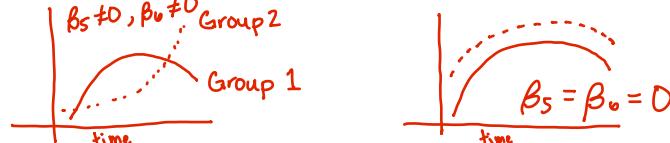
Using polynomial trends such as linear or quadratic, we can model longitudinal data as a function of time. Linear trends are the most common and interpretable ways to model change in mean over time. In an example comparing a treatment group to a control group, we can fit a linear trend using the following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

*I think you  
have some  
confusion  
here*

Unclear ...  
 $\beta_4 = 0 \Rightarrow$   
any difference in response  
btwn group is constant  
over time.

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.

### 1.3.5 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

Use notation  
 to specify distribs  
 of  $b_i$  and  $\epsilon_{ij}$ .  
 Clarify each are  
 iid and  $b_i$   
 are indep. of  $\epsilon_{ij}$ .  
 Incorporate interpretation of  
 $b_i$  in same paragraph.

is unique to each individual.  $b_i$  is interpreted as how the subject <sup>is intercept</sup> 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

$$\text{Cov}(Y_{ij}) = \begin{pmatrix} \sigma_b^2 + \sigma^2 & \sigma_b^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma^2 & \sigma_b^2 + \sigma^2 & \sigma_b^2 & \dots & \sigma_b^2 \\ \sigma_b^2 & \sigma_b^2 & \sigma_b^2 + \sigma^2 & \dots & \dots \end{pmatrix}$$

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

$$Y_i = X_i\beta + Z_i b_i + \epsilon_i$$

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

$$\Sigma_i = \text{Cov}(Y_i) = \text{Cov}(Z_i b_i) + \text{Cov}(\epsilon_i) = Z_i G Z_i' + R_i.$$

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

$$Y_i = X_i \beta + Z_i b_i + \epsilon_i,$$

where both  $X_i$  and

$$X_i = Z_i = \begin{pmatrix} 1 & t_{i1} \\ 1 & t_{i2} \\ \dots & \dots \\ 1 & t_{in_i} \end{pmatrix}$$

so the equation

For the  $i^{th}$  subject at the  $j^{th}$  measurement, the equation is as follows:

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + b_{1i} + b_{2i} t_{ij} + \epsilon_{ij}.$$

Convention is to start with  $\beta_0$  and  $b_0$  } recognizable as intercept terms

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 com-

ponents 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.

More to  
Ch 2

### estimation section

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:

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

don't need  
parents; i should  
be attached to  $\Sigma$

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:

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

Cite FLW  
Somewhere  
in here.

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.

Can be  
deleted  
(no extension  
is shown so  
can't follow  
description)

## **1.4 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.

## **1.5 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.