



VCU

Virginia Commonwealth University
VCU Scholars Compass

Theses and Dissertations

Graduate School

2008

Tolerance Intervals in Random-Effects Models

Kakotan Sanogo
Virginia Commonwealth University

Follow this and additional works at: <https://scholarscompass.vcu.edu/etd>



Part of the [Biostatistics Commons](#)

© The Author

Downloaded from

<https://scholarscompass.vcu.edu/etd/1661>

This Thesis is brought to you for free and open access by the Graduate School at VCU Scholars Compass. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of VCU Scholars Compass. For more information, please contact libcompass@vcu.edu.

TOLERANCE INTERVALS IN RANDOM-EFFECTS MODELS

A Thesis submitted in partial fulfillment of the requirements for the degree of Master of
Science at Virginia Commonwealth University.

by

KAKOTAN SANOGO

B.S. Computer Engineering, Minnesota State University, 2004

Director: JESSICA M. KETCHUM, PH.D.

ASSISTANT PROFESSOR, DEPARTMENT OF BIostatISTICS

Virginia Commonwealth University

Richmond, Virginia

December 2008

Acknowledgement

Above all I would like to thank GOD almighty for allowing me to live up to this time, in good mental and physical health. I dearly thank my lovely parents for their many blessings and the principles they've taught me for many years. To me and my four younger siblings they've been wonderful educators. I cannot thank them enough for sending me to America to pursue my dreams. I would like to thank my cousin and longtime friend and mentor, Dr. Oumar Sy, for talking to me about the fascinating field of Biostatistics during my short stop at his house in the summer of 2004. Without his enlightenment, this work would not have taken place. I would like to greatly thank my advisor, Dr. Jessica M. Ketchum, for taking time of her super busy schedule to assist me throughout this thesis. Without her dedication and kindness, this work would be but a dream unrealized. I would like to equally thank Dr. Charles (Charlie) W. Kish for bringing up the topic of this thesis to my attention and for being my mentor during the course of my year-long internship at Wyeth Consumer Healthcare. Without his guidance, I would probably still be looking for a research topic. I would also like to thank Dr. Ramakrishnan Wiswanathan (Ramesh) for providing valuable advise throughout the course of this thesis and for being a valuable source of knowledge through the advanced applied statistical courses I took from him during my training as a graduate student. I would like to equally thank all the professors of the department of Biostatistics for sharing their invaluable knowledge and for always keeping their door open. I would also like to thank all my fellow colleagues and the department secretaries, for their assistance and kindness was invaluable to my survival

during my journey as a graduate student. I would like to thank my many friends that understood my silence during this defining moment of my life. I would finally like to thank my fiancée Fatoumata (Fatim) for being kind and patient during this long journey. To you all I owe much and bow to your wisdom.

Table of Contents

	Page
Acknowledgements	ii
List of Tables.....	vi
List of Figures	vii
Abstract	viii
1 Introduction	1
1.1. Introduction to Statistical Intervals	1
1.2. Motivation	2
1.3. Objective	6
2 Methodology	8
2.1. Introduction	8
2.2. The Normal Fixed-Effects Model Definition.....	8
2.3. The Normal Random-Effects Model Definition	9
2.4. Confidence Interval Definition	10
2.5. Prediction Interval Definition	10
2.6. Tolerance Interval	11
2.6.1 Introduction	11
2.6.2 Wilks' Method for Tolerance Interval.....	12
2.6.3 Graybill's Method for Tolerance Interval.....	13
2.6.4 Jonsson's Method for Tolerance Interval	13
2.6.4.1 Introduction.....	13

2.6.4.2	Estimation	14
2.6.3.2.1	Estimation of \hat{V}	17
2.6.3.2.2	Estimation of \hat{K}_L	19
2.6.3.3	β -Expectation Tolerance Interval	21
3	Applications	22
3.1.	Presenting Results from the Analysis	23
3.2.	Discussion and summary of the Analysis	28
4	Conclusion and Future Research.....	30
5	List of References	33
6	Appendix	35
6.1.	SAS Code.....	35

List of Tables

	Page
Table 1: Data Structure	22
Table 2: 95% Confidence Intervals.....	26
Table 3: 95% Tolerance Intervals	26

List of Figures

Page

Figure 1: Tolerance Intervals27

Abstract

INTERVALS ESTIMATES IN RANDOM-EFFECTS MODELS

A Thesis submitted in partial fulfillment of the requirements for the degree of Master of Science at Virginia Commonwealth University.

Virginia Commonwealth University, 2008

Major Director

Jessica M. Ketchum

Assistant Professor Department of Biostatistics

In the pharmaceutical setting, it is often necessary to establish the shelf life of a drug product and sometimes suitable to assess the risk of product failure at the desired expiry period. The current statistical methodology use confidence intervals for the predicted mean to establish the expiry period and prediction intervals for a predicted new assay value or a tolerance interval for a proportion of the population for use in a risk assessment. A major concern is that most methodology treat a homogeneous subpopulation, say batch, either as a fixed effect and therefore uses a fixed-effects regression model (Graybill, 1976) or as a mixed-effects model limited to balanced data

structures (Jonsson, 2003). However, batch is definitely a random effect as this fact has been reflected by some recent methodology [Altan, Cabrera and Shoung (2005), Hoffman and Kringle (2005)]. Thus, to assess the risk of product failure at expiry, it is necessary to use tolerance intervals since they provide an estimate of the proportion of assay values and/or batches failing at the expiry period. In this thesis, we illustrate the methodology described by Jonsson (2003) to construct β -expectation tolerance limits for longitudinal data in a random-effects setting. We underline the limitations of Jonsson's approach to constructing tolerance intervals and highlight the need for a better methodology.

1 Introduction

1.1. Introduction to Interval Estimates

Scientists and engineers frequently express the need to quantify the uncertainty associated with a point estimate in order to make decisions from limited sample data. They may wish to obtain more data prior to making a decision if their knowledge of the uncertainty is imprecise. To quantify such uncertainty, interval estimates are constructed around a point estimate. Three different types of interval estimates may be calculated from sampled data. Depending upon the type of application, the analyst may choose a confidence interval, a tolerance interval, or a prediction interval.

Using sample data, some researchers may be interested in estimating a confidence interval, a range of values expected to encompass the population parameter of interest with some specified level of confidence. One way to think about confidence intervals is to consider drawing many samples (in the same manner) from a population. Each sample yields its own estimate of the parameter of interest (e.g., the population mean) and corresponding confidence interval with a selected or desired confidence coefficient (e.g., 95%). In this repeated sense, approximately 95% of the confidence intervals will enclose the population mean.

Similar to the confidence interval, a tolerance interval is a range of values expected to contain a certain percentage of observations from a population on the average. For example, one may be interested in determining a range of values expected to encompass 90% of the population on the average.

Confidence intervals and tolerance intervals are both interval estimates for parameters of the population. A prediction interval is a range of values expected to encompass a new (future) observation from the population with a specified level of confidence. For example, one may be interested in determining a range of values containing the next predicted value with 95% confidence.

In order to choose the most appropriate interval (confidence, tolerance, or prediction), the analyst must decide whether the main interest of the application resides in describing the population from which the sample has been selected or in predicting the results of a future sample from the same population.

1.2. Motivation

In the pharmaceutical setting, it is often necessary to establish the expiry period of a drug and sometimes suitable to assess the risk of product failure at the desired expiry period. The Food and Drug Administration (FDA) guidelines (FDA, 1987) requires that a minimum of three batches be tested in stability analysis to account for batch-to-batch variability so that a single shelf life is applicable to all future batches manufactured under similar circumstances. In addition to the estimation of the individual shelf life for each batch, it is also desirable to establish a single shelf life for a drug product based on combined stability data from all batches. The FDA guidelines requires that preliminary tests of batch similarity be performed before combining the stability data from all batches. A test for differences in the intercepts and differences in the slopes of degradation lines among different batches is performed to evaluate batch similarity. The

FDA recommends the 0.25 level of significance to test these hypotheses. Thus, the single shelf life can be determined, based on the ordinary least-squares methods, as the time point at which the 95% lower confidence bound for the mean degradation curve of the drug characteristic intersects the approved lower specification limit.

If the hypotheses of equal intercepts and equal slopes are not rejected at the 0.25 level of significance, a single expiration dating period is usually estimated by fitting a single degradation curve based on the pooled stability data of all batches under the assumption that batch effects are fixed. If the hypotheses of equal intercepts and equal slopes are rejected at the 0.25 level of significance, the FDA recommends determining a single expiration dating period of the drug product based on the minimum of shelf lives obtained from each batch. However, Chow and Shao (1991) showed that this method had no statistical justification since the minimum approach is conservative and does not take into account batch-to-batch variability.

Confidence intervals for the predicted mean are commonly used to establish the expiry period of a drug product. Prediction intervals for a predicted new assay value and tolerance intervals for a proportion of the population are sometimes used in risk assessment. To assess the risk of product failure at expiry, it is more appropriate to use tolerance intervals since they provide an interval estimate for the proportion of assay values in the population failing at the expiry period.

The need for tolerance intervals was greatly emphasized during the first half of the twentieth century. Wilks (1942) defined and constructed tolerance limits in the case of normal distribution with unknown mean and variance. The use of tolerance intervals

based on linear models became the interest of various researchers such as Wallis (1951), Weissberg and Beatty (1960), Lieberman and Miller (1963), Ellison (1964), Howe (1969), Graybill (1976). A procedure for establishing a two-sided tolerance interval based on a balanced mixed-effects model was proposed by Liao and Iyer (2004). Under the assumption of fixed batch effects, Hsu and Ruberg (1992) proposed a method to estimate the expiration dating period of a drug product by using multiple comparison technique for pooling stability data with the worst batch. The foregoing methodology (fixed batch effects model) assumes that the drug characteristic decreases linearly over time. The comparison of regression lines necessitates not only a test of equality of intercepts and equality of slopes but also the equivalence of within batch variability. It should be recognized, however, that the between-batch variation is often ignored during the decision making process for pooling stability data across batches.

The FDA guidelines indicate that the batches used in long-term stability studies for the establishment of drug shelf life should constitute a random sample from the population of future production batches. The FDA also requires that all estimated expiration dating periods should be applicable to all future batches. Under these assumptions, the statistical methods derived from the fixed-effects models may not be appropriate. This is due to the fact that statistical inferences about the expiration dating period obtained from a fixed-effects model can only be drawn from the batches under study and cannot necessarily be applied to future or unobserved batches. The use of statistical methods based on a random-effects model is therefore more appropriate for establishing the expiration dating period for future production batches. Several

researchers have gained interest in the use of tolerance intervals in random-effects settings. Lemon (1977) and Mee and Owen (1983) considered the case of one-sided tolerance intervals for balanced one-way random-effects models. This was soon followed by an extension to the unbalanced random-effects model [see Bhaumick and Kulkarni (1991,1996), Bagui et al. (1996)]. Two-sided tolerance intervals for balanced one-way random-effects models were also considered by Mee (1984) and an extension for unbalanced data was described by Beckman and Tietjen (1989) and Wang and Iyer (1994). The computation of a one-sided tolerance limit for a one-way random-effects model for both balanced and unbalanced data using the concept of a generalized confidence interval explored by Weerahandi (1993, 1995) was extended by Krishnamoorthy and Thomas (2004). Hoffman and Kringle (2005) proposed a methodology for constructing two-sided tolerance intervals for general random-effects models in both balanced and unbalanced cases. A procedure for constructing a two-sided tolerance limits without the normality assumption for both balanced and unbalanced ANOVA models by using a nested bootstrap method was proposed by Shoung et al. (2005). However, all the aforementioned analytical methods are based on a cross-sectional approach and therefore do not utilize the longitudinal structure of the data (which can lead to more accurate tolerance intervals) or used distribution-free methods which have limitations in small samples since they are based on order statistics. The last two aforementioned methodologies [Hoffman and Kringle (2005), and Shoung et al. (2005)] use the β -content tolerance intervals procedure which are mainly intended for drugs where the risk of adverse side effects rapidly increases with an overdose, i.e. even a

minor overdose may result in death (Petzold, 2001). Jonsson (2003) proposed a new methodology that not only took into account the longitudinal structure of the data but also used the β -expectation tolerance intervals procedure which are intended for drugs where the expected outcome of an overdose may not cause death (Petzold, 2001). However, Jonsson's approach to tolerance intervals treats the slopes as a fixed effect and needs to be enhanced since the random-effects slopes are of extreme importance to the pharmaceutical industry.

In summary, the main difference between the fixed-effects models and the random-effects models is that the random-effects model incorporates the fact that batches are considered a random sample drawn from the population of all production batches, including future ones if the process does not change. Hence, the intercepts and slopes used to characterize the degradation of a drug product should be considered as random variables.

1.3. Objective

The aim of the present work is to describe and illustrate tolerance interval methods based on random-effects and fixed-effects models. We will also describe methods for confidence intervals and prediction intervals based on fixed-effects models. Methods for tolerance intervals in the random-effects setting will be based on those described by Jonsson (2003) while models for tolerance intervals in the fixed-effects setting will be based on those described by Wilks (1941) and Graybill (1976). In Chapter 2, we will introduce the fixed-effects model, the random-effects model, and the interval

estimates. In Chapter 3, we will demonstrate the results from the methodology using an analysis dataset. In Chapter 4, we will draw the conclusion from the analysis and make suggestions for future research.

2 Methodology

2.1. Introduction

In this chapter, we introduce the fixed-effects model (Section 2.2) and the random-effects model (Section 2.3). Methods for describing confidence intervals, prediction intervals, and tolerance intervals are described in Sections 2.4, 2.5, and 2.6 respectively.

2.2. The Normal Fixed-Effects Model Definition

The fixed-effects model is of the form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}, \quad (2.1)$$

where

\mathbf{y} is an $n \times 1$ vector of observed response values,

\mathbf{X} is an $n \times p$ ($n > p$) observed design matrix corresponding to the fixed-effects,

$\boldsymbol{\alpha}$ is a $p \times 1$ vector of fixed-effects parameters, and

$\boldsymbol{\varepsilon}$ is an $n \times 1$ unobservable vector of residuals.

The residuals $\varepsilon_i, i = 1, \dots, n$ are assumed to be independent and identically

distributed $N(0, \sigma^2)$. Thus, the variance of \mathbf{y} , $\text{var}(\mathbf{y}) = \mathbf{V}$ is given by

$\mathbf{V} = \text{var}(\mathbf{X}\boldsymbol{\alpha} + \boldsymbol{\varepsilon}) = \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$, where \mathbf{I} is an $n \times n$ identity matrix.

2.3. The Normal Random-Effects Model Definition

The random-effects model introduced by Laird and Ware (1982) extends the fixed-effects model in equation (2.1) such that

$$\mathbf{y} = \mathbf{X}\boldsymbol{\alpha} + \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (2.2)$$

where

y , X , $\boldsymbol{\alpha}$, and $\boldsymbol{\varepsilon}$ are as defined in the fixed effects model,

Z is an $n \times q$ observed design matrix for the random-effects, and

$\boldsymbol{\beta}$ is the $q \times 1$ vector of random-effects/coefficients parameters.

The covariance of \mathbf{y} , $\text{var}(\mathbf{y}) = \mathbf{V}$ given by $\mathbf{V} = \text{var}(\mathbf{X}\boldsymbol{\alpha} + \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\varepsilon})$. This model assumes that the random-effects and the residuals are independent so that

$$\mathbf{V} = \text{var}(\mathbf{X}\boldsymbol{\alpha}) + \text{var}(\mathbf{Z}\boldsymbol{\beta}) + \text{var}(\boldsymbol{\varepsilon}).$$

Since $\boldsymbol{\alpha}$ describes the fixed-effects parameters, $\text{var}(\mathbf{X}\boldsymbol{\alpha}) = 0$. Hence,

$$\mathbf{V} = \mathbf{Z}\text{var}(\boldsymbol{\beta})\mathbf{Z}' + \text{var}(\boldsymbol{\varepsilon}).$$

Under the assumption that the random-effects follow normal distributions and letting

$\text{var}(\boldsymbol{\beta}) = \mathbf{G}$ we obtain

$$\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}' + \boldsymbol{\Sigma}, \quad (2.3)$$

where $\boldsymbol{\beta} \sim \mathbf{N}(\mathbf{0}, \mathbf{G})$ and $\boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \boldsymbol{\Sigma})$.

2.4. Confidence Interval Definition

For a fixed-effects simple linear regression model, consider estimating the mean response of a population given particular values of the predictor x . A two-sided $100(1 - \alpha)\%$ confidence interval for the mean response $\mu_{y|x_0}$ is given by

$$\hat{y}_0 \pm t_{1-\frac{\alpha}{2}} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \quad (2.4)$$

where \hat{y}_0

represents the estimate of the mean response at $x = x_0$,

n is the population size, and

$t_{1-\frac{\alpha}{2}}$ is the $1-\frac{\alpha}{2}$ percentile of the central t-distribution with $n-2$ degrees of freedom.

2.5. Prediction Interval Definition

For the fixed-effects simple linear regression model, consider predicting the response of a single future observation given particular values of the predictor x . A two-sided $100(1 - \alpha)\%$ prediction interval for a single response value Y_0 is given by

$$\hat{y}_0 \pm t_{1-\frac{\alpha}{2}} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \quad (2.5)$$

where

\hat{y}_0 represents the estimate of the mean response at $x = x_0$,

n is the population size, and

$t_{1-\frac{\alpha}{2}}$ is the $1-\frac{\alpha}{2}$ percentile of the central t-distribution with $n-2$ degrees of freedom.

Confidence intervals reflect the uncertainty in estimating the mean response which is a function of the parameter estimates. Prediction intervals reflect the uncertainty from the parameter estimates as well as the uncertainty of a future observation. This is why prediction intervals are wider than confidence intervals. Comparing the two formulae in equation (2.4) and equation (2.5), the prediction intervals include an extra “1” in the square root which reflects the additional uncertainty.

2.6. Tolerance Interval

2.6.1 Introduction

In this section, we first define the tolerance interval introduced by Wilks (1941) as it applies to cross-sectional data under the fixed-effects model. We then define the β -content tolerance interval described by Graybill (1976) under the fixed-effects model as it applies to longitudinal data. We finally describe the methods used for estimation of the β -expectation tolerance interval as it applies to longitudinal data in the normal random-effects setting as it was introduced by Jonsson (2003).

Consider estimating a number γ_p such that a $(1-p)$ proportion of the responses in the population under study is below it; or two numbers $\gamma_{1-p/2}, \gamma_{p/2}$ such that a $(1-p)$ proportion of the responses in the population is between the two numbers. Let the assay result (e.g., percent of label claim) be represented by a continuous random

variable Y with cumulative distribution function F , and let w_1 and w_2 be defined by the relation $P_Y(w_1 < Y < w_2) = F(w_1) - F(w_2) = \beta$, where β is a predetermined probability. A β -expectation tolerance interval requires that $100\beta\%$ of the individual responses from the batches fall between the estimated limits on the average; that is, the expectation over \hat{w}_1 and \hat{w}_2 , $E(P_Y(\hat{w}_1 < Y < \hat{w}_2))$, equals β exactly or at least approximately (Jonsson, 2003).

Though β -expectation tolerance intervals can be constructed with or without distributional assumptions, distribution-free tolerance limits have limitations in small samples, as they are based on order statistics $(Y_{(r)}, Y_{(n-r+1)})$ with $r < (n+1)/2$.

2.6.2 Wilks' Method for Tolerance Interval

Let the independent and identically distributed random variables Y_1, \dots, Y_n from a normal distribution $N(\mu, \sigma^2)$ represent a cross-sectional sample at some time. If Y is a new observation from $N(\mu, \sigma^2)$ and if \bar{Y} and S^2 are the unbiased estimators of μ and σ^2 ,

respectively, then $\frac{Y - \bar{Y}}{S}$ is distributed as $\sqrt{1 + \frac{1}{n}} T_{(n-1)}$, where $T_{(n-1)}$ denotes a Student T-

variable with $n-1$ degrees of freedom. Hence, the β -expectation tolerance bounds at each time point introduced by Wilks (1941) are of the form

$$\bar{Y} \pm K_C S, \quad (2.6)$$

where

$$K_C = \sqrt{1 + \frac{1}{n}} t_{\frac{1+\beta}{2}, (n-1)} \quad (2.7)$$

and $t_{\frac{1+\beta}{2}(n-1)}$ is the $100(\frac{1+\beta}{2})$ -percentile in the distribution of $T_{(n-1)}$.

2.6.3 Graybill's Method for Tolerance Interval

For a fixed-effects simple linear regression model, consider finding a range of values expected to contain a certain percentage of observations from a population with some specified level of confidence. The Q^{th} tolerance limits at the point x_0 and with confidence coefficient $1 - \alpha$ is of the form

$$\hat{\alpha}_0 + \hat{\alpha}_1 x_0 \pm g_Q \hat{\sigma} \quad (2.8)$$

where g_Q is given by

$$g_Q = A t_{\alpha; n-p; (\theta)},$$

$$A^2 = \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum (x_i - \bar{x})^2},$$

$t_{\alpha; n-p; (\theta)}$ is the upper α probability point of the noncentral t distribution with $n - p$ degrees

of freedom and noncentrality parameter $\theta = \frac{N_Q}{A}$.

2.6.4 Jonsson's Method for Tolerance Interval

2.6.4.1 Introduction

A considerable improvement of Wilks' method for tolerance intervals can be achieved by utilizing the longitudinal structure of the data so that all the n subjects are

used at each time point. Using the pharmaceutical stability framework, a mixed-effects model, the variance components regression model, used by Jonsson (2003) is defined as

$$Y_{ij} = \mu_j + \delta t_i + \varepsilon_{ij}, \quad i = 1, \dots, T, \quad j = 1, \dots, n, \quad (2.9)$$

where

Y_{ij} is the assay result (percent label claim) for the j th batch at the i th time point,

μ_j is a random-effect that reflects factors that are specific to the j th batch,

δ is a fixed-effect which expresses the change over time common to all batches,

t_i represent the time points at which the response is defined, and

ε_{ij} is a random residual that summarizes the effects of all factors that have not been

included in the model.

The β -expectation tolerance bounds at time t for the random-effects model are given by

$$\hat{\mu}_t \pm \hat{K}_L \hat{V}, \quad (2.10)$$

where $\hat{\mu}_t$ and \hat{V} are estimators of the mean and the square root of the variance of Y_{ij} ,

respectively, and are based on all nT observations. The quantity \hat{K}_L depends on the time t

and is estimated from the data whereas K_C in equation (2.7) is a constant across

time. \hat{K}_L is determined so that the β -expectation property holds.

2.6.4.2 Estimation

The problem with constructing tolerance bounds using equation (2.9) are two fold: to find a good estimator of V and to determine the value of \hat{K}_L , which depends upon the distribution of $\frac{Y_{ij} - \hat{\mu}_t}{\hat{V}}$. One approach is to use a Taylor series approximation so that the tolerance bounds have at least approximately the β -expectation property.

We assume that the $\mu_j, j = 1, \dots, n$, are identically distributed $N(\alpha, \sigma_U^2)$ and the $\varepsilon_{ij}, j = 1, \dots, n, i = 1, \dots, T$, are identically distributed $N(0, \sigma_\varepsilon^2)$. We also assume that the μ_j and the ε_{ij} are independent of each other. Thus the Y_{ij} are $N(\alpha + \delta t_i, \sigma_U^2 + \sigma_\varepsilon^2)$. The tolerance interval is estimated as $(\hat{w}_1, \hat{w}_2) = (\hat{\mu}_i - \hat{K}_L \hat{V}, \hat{\mu}_i + \hat{K}_L \hat{V})$, where $\hat{\mu}_i = \hat{\alpha} + \hat{\delta} t_i$ is an estimator of $\mu_i = \alpha + \delta t_i$ and \hat{V} is some unbiased estimator of $V = \sqrt{\sigma_U^2 + \sigma_\varepsilon^2}$.

Based on data from the j^{th} batch, let $\hat{\alpha}_j$ and $\hat{\delta}_j$ be the ordinary least squares estimators of α and δ . The best unbiased estimators of α and δ , $\hat{\alpha}$ and $\hat{\delta}$, are given by $\hat{\alpha} = \sum_{j=1}^n \frac{\hat{\alpha}_j}{n}$ and

$\hat{\delta} = \sum_{j=1}^n \frac{\hat{\delta}_j}{n}$, respectively. In addition, we define the corrected sum of squares as

$W_{tt} = \sum_{i=1}^T (t_i - \bar{t})^2$, where $\bar{t} = \sum_{i=1}^T \frac{t_i}{T}$ is the mean time. Then,

$$\hat{\alpha} + \hat{\delta} t_i \sim N\left(\alpha + \delta t_i, \frac{\sigma_U^2 + \sigma_\varepsilon^2}{n} + \frac{\sigma_\varepsilon^2 (t_i - \bar{t})^2}{n W_{tt}}\right).$$

Let $\hat{w}'_r = \mu_i \pm K_L \hat{V}$, ($r = 1, 2$) be provisory estimators such that K_L is an arbitrary constant. Since $V(\hat{w}'_r) = V(\hat{w}')$, using a Taylor series approximation we obtain

$$E\left(F(\hat{w}'_2) - F(\hat{w}'_1)\right) \approx F(w_2) - F(w_1) + \frac{V(\hat{w}')}{2} \left((F^{(2)}(w_2) - F^{(2)}(w_1)) \right). \quad (2.11)$$

If we let F be the cdf of a normally distributed random variable with variance V^2 and Φ be the cdf of a corresponding standardized variable, we obtain the following three relations:

$$F(w_2) - F(w_1) = 2\Phi(K_L) - 1,$$

$$F^{(2)}(w_2) - F^{(2)}(w_1) = \frac{-2K_L \exp\left\{-\frac{K_L^2}{2}\right\}}{V^2 \sqrt{2\pi}}, \text{ and}$$

$$\begin{aligned} V(\hat{w}') &= V(\mu_i \pm K_L \hat{V}) = V(\hat{\alpha} + \hat{\delta}t_i \pm K_L \hat{V}) \\ &= V(\hat{\alpha} + \hat{\delta}t_i) + K_L^2 V(\hat{V}) \\ &= \frac{\sigma_u^2 + \sigma_\varepsilon^2}{n} + \frac{\sigma_\varepsilon^2 (t_i - \bar{t})^2}{nW_n} + K_L^2 V(\hat{V}). \end{aligned}$$

Thus,

$$\begin{aligned}
E(F(\hat{w}'_2) - F(\hat{w}'_1)) &= 2\Phi(K_L) - 1 - \frac{1}{\sqrt{2\pi}} \left[\frac{\sigma_U^2 + \sigma_\varepsilon^2}{nV^2} + \frac{\sigma_\varepsilon^2(t_i - \bar{t})^2}{nV^2W_{tt}} + \frac{K_L^2V(\hat{V})}{V^2} \right] K_L \exp\left\{-\frac{K_L^2}{2}\right\} \\
&= 2\Phi(K_L) - 1 - \frac{1}{\sqrt{2\pi}} \left[\frac{1}{n} + \frac{\sigma_\varepsilon^2}{V^2} \frac{(t_i - \bar{t})^2}{nW_{tt}} + \frac{K_L^2V(\hat{V})}{V^2} \right] K_L \exp\left\{-\frac{K_L^2}{2}\right\} \\
&\approx 2\Phi(K_L) - 1 + \frac{1}{2} \left[\frac{\sigma_U^2 + \sigma_\varepsilon^2}{n} + \frac{\sigma_\varepsilon^2(t_i - \bar{t})^2}{nW_{tt}} + K_L^2V(\hat{V}) \right] \left[\frac{-2K_L \exp\left\{-\frac{K_L^2}{2}\right\}}{V^2\sqrt{2\pi}} \right] \\
&= 2\Phi(K_L) - 1 - \frac{1}{\sqrt{2\pi}} \left[\frac{\sigma_U^2 + \sigma_\varepsilon^2}{nV^2} + \frac{\sigma_\varepsilon^2(t_i - \bar{t})^2}{nV^2W_{tt}} + \frac{K_L^2V(\hat{V})}{V^2} \right] K_L \exp\left\{-\frac{K_L^2}{2}\right\} \\
&= 2\Phi(K_L) - 1 - \frac{1}{\sqrt{2\pi}} \left[\frac{1}{n} + \frac{\sigma_\varepsilon^2}{V^2} \frac{(t_i - \bar{t})^2}{nW_{tt}} + K_L^2 \frac{V(\hat{V})}{V^2} \right] K_L \exp\left\{-\frac{K_L^2}{2}\right\}.
\end{aligned} \tag{2.12}$$

2.6.4.2.1. Estimation of \hat{V}

We define the following sums of squares:

$$W_{YY} = \sum_{j=1}^n \sum_{i=1}^T (Y_{ij} - \bar{Y}_j)^2, \tag{2.13}$$

$$W_{tY} = \sum_{j=1}^n \sum_{i=1}^T (t_i - \bar{t})(Y_{ij} - \bar{Y}_j), \quad \text{and} \tag{2.14}$$

$$S = \sum_{j=1}^n (\bar{Y}_j - \bar{Y})^2. \tag{2.15}$$

If we define $\bar{Y} = \sum_{j=1}^n \frac{\bar{Y}_j}{n}$, then, $\frac{S}{n-1} \sim (\sigma_U^2 + \frac{\sigma_\varepsilon^2}{T}) \frac{\chi_{n-1}^2}{n-1}$, and.

$$\hat{\sigma}_\varepsilon^2 = \frac{W_{YY} - \delta W_{tY}}{n(T-1) - 1} \sim \frac{\sigma_\varepsilon^2}{n(T-1) - 1} \chi_{n(T-1)-1}^2. \tag{2.16}$$

Since $\frac{S}{n-1}$ and $\hat{\sigma}_\varepsilon^2$ are independent of each other and of $\hat{\alpha}$ and $\hat{\delta}$, an unbiased estimator

of $V^2 = \sigma_U^2 + \sigma_\varepsilon^2$ is given by

$$\hat{V}^2 = \hat{\sigma}_U^2 + \hat{\sigma}_\varepsilon^2 = \frac{S}{n-1} + \hat{\sigma}_\varepsilon^2 \left(1 - \frac{1}{T}\right). \quad (2.17)$$

Jonsson (2003, Table 1) shows that $\sqrt{\hat{V}^2}$ is in fact a biased estimator of V and therefore needs to be corrected. Using a Taylor series expansion, the expectation of V is

$$E(\hat{V}) = E(\sqrt{\hat{V}^2}) \approx \sqrt{V^2} - \frac{V(\hat{V}^2)}{8(V^2)^{3/2}} = V \left[1 - \frac{V(\hat{V}^2)}{8(V^2)^2}\right]. \quad (2.18)$$

This leads to the following adjusted estimator of V

$$\sqrt{\hat{V}^2} \left[1 - \frac{V(\hat{V}^2)}{8(V^2)^2}\right]^{-1}.$$

It follows from equations (2.16) and (2.17) that the variance of \hat{V}^2 , $V(\hat{V}^2)$, is of the form

$$\begin{aligned} V(\hat{V}^2) &= V\left(\frac{S}{n-1} + \hat{\sigma}_\varepsilon^2 \left(1 - \frac{1}{T}\right)\right) \\ &= V\left(\frac{\sigma_U^2 + \frac{\sigma_\varepsilon^2}{T}}{n-1} + \frac{\sigma_\varepsilon^2 \left(1 - \frac{1}{T}\right)}{n(T-1)-1}\right) \\ &= 2\left[\frac{(\sigma_U^2 + \frac{\sigma_\varepsilon^2}{T})^2}{n-1} + \frac{\sigma_\varepsilon^4 \left(1 - \frac{1}{T}\right)^2}{n(T-1)-1}\right]. \end{aligned}$$

This implies that

$$\begin{aligned}
\frac{V(\hat{V}^2)}{8(V^2)^2} &= \frac{V(\hat{V}^2)}{8(\sigma_U^2 + \sigma_\varepsilon^2)^2} \\
&= \frac{2}{8} \left[\frac{(\sigma_U^2 + \frac{\sigma_\varepsilon^2}{T})^2}{(n-1)(\sigma_U^2 + \sigma_\varepsilon^2)^2} + \frac{\sigma_\varepsilon^4(1 - \frac{1}{T})^2}{[n(T-1)-1](\sigma_U^2 + \sigma_\varepsilon^2)^2} \right] \\
&= \frac{1}{4} \left[\frac{(1-R)^2}{n-1} + \frac{R^2}{n(T-1)-1} \right],
\end{aligned}$$

where

$$R = \frac{\sigma_\varepsilon^2(1 - \frac{1}{T})}{(\sigma_U^2 + \sigma_\varepsilon^2)}. \quad (2.19)$$

Hence an estimator of $V = \sqrt{\sigma_U^2 + \sigma_\varepsilon^2}$ is given by

$$\hat{V} = \sqrt{\hat{\sigma}_U^2 + \hat{\sigma}_\varepsilon^2} \left[1 - \frac{1}{4} \left\{ \frac{(1-\hat{R})^2}{n-1} + \frac{\hat{R}^2}{n(T-1)-1} \right\} \right]^{-1}, \quad (2.20)$$

where

$$\hat{R} = \frac{\hat{\sigma}_\varepsilon^2(1 - \frac{1}{T})}{\hat{\sigma}_U^2 + \hat{\sigma}_\varepsilon^2} \cdot \left(\frac{n(T-1)-3}{n(T-1)-1} \right). \quad (2.21)$$

This correction factor involving the estimator \hat{R} helps reduce the bias in small samples since the ratio of the estimated variance components has an F-distribution. Jonsson (2003) found through simulation studies that the uncorrected estimator of V has a negative bias that decreases in absolute value with increasing values of n , and to a lesser extent with increasing values of T .

2.6.4.2.2. Estimation of \hat{K}_L

Using equation (2.12) and equating it to a predetermined value of β , we may get a relationship between R and K_L .

$$\begin{aligned}
 E(F(\hat{w}_2) - F(\hat{w}_1)) &= 2\Phi(K_L) - 1 - \\
 &\frac{1}{\sqrt{2\pi}} \left[\frac{1}{n} \left\{ 1 + \frac{R}{1 - \frac{1}{T}} + \frac{(t - \bar{t})^2}{W_u} \right\} + \frac{K_L^2}{2} \left\{ \frac{(1-R)^2}{n-1} + \frac{R^2}{n(T-1)-1} \right\} \right] + K_L \exp\left\{-\frac{K_L^2}{2}\right\} \\
 &= \beta
 \end{aligned} \tag{2.22}$$

where R is given by (2.19).

We then use this relationship to find an estimated value of K_L from an observed value of R . We write a quadratic equation in R by equating (2.12) to the predetermined value β since it is difficult to write K_L as an explicit function of R . This quadratic equation has one root of interest.

Let

$$C = \frac{(t_i - \bar{t})^2}{n(1 - \frac{1}{T})W_u}, \text{ and} \tag{2.23}$$

$$A = \frac{2}{K_L^2} \left[\frac{\sqrt{2\pi}(2\Phi(K_L) - 1 - \beta)}{K_L \exp\{-\frac{K_L^2}{2}\}} - \frac{1}{n} \right], \tag{2.24}$$

where C is a component that is determined by spacing of the time points at which the measurements are made. For a given n and T , $C=0$ at $t = \bar{t}$ and reaches the maximum at the end points of the range of t . (See Jonsson, 2003 for an example). We could note,

however, that the maximum at the end points decreases with increasing T (Jonsson, 2003). Equation (2.22) becomes

$$R = \frac{(n(T-1)-1)\left[1 - \frac{(n-1)C}{K_L^2}\right] - [(n(T-1)-1)^2\left[1 - \frac{(n-1)C}{K_L^2}\right]^2 - (nT-2)(n(T-1)-1)(1-(n-1)A)]^{\frac{1}{2}}}{nT-2}. \quad (2.25)$$

There exists a one-to-one relationship between $R \in (0,1)$ and $K_L \in I_{K_L}$ such that R decreases with increasing K_L . The quantity \hat{K}_L denotes the value of K_L corresponding to an estimated value of R given in (2.21). Jonsson (2003) states that K_L becomes empirically almost linearly dependent on

$$Z = \frac{(1-R)^2}{n-1} + \frac{R^2}{n(T-1)-1} \quad (2.26)$$

provided the predetermined value (β) is 0.90 or 0.95. For larger values of β , especially for small n , we fit K_L to polynomial functions of Z . The estimation of K_L is done using a computer program which is found in the Appendix section.

2.6.4.3 β -Expectation Tolerance Interval

The estimated 95%-expectation tolerance interval in the population of random batches at time t based on longitudinal data is of the form $(\hat{\alpha} + \hat{\delta}_t - \hat{K}_L \hat{V}, \hat{\alpha} + \hat{\delta}_t + \hat{K}_L \hat{V})$.

The values of \hat{V} and \hat{K}_L are derived from equations (2.20) and (2.25) respectively.

3 Applications

In this chapter, the methodology for estimation proposed in Chapter 2 is studied using an example dataset. We demonstrate the results using a random-effects model with data from a pharmaceutical stability data simulation which was performed by Obenchain (1990).

Table 1: Data Structure

Observations	Batch	Month	\bar{Y}
1	1	0	102.783
2	1	1	99.350
3	1	3	98.625
4	1	6	101.525
5	1	9	96.750
6	1	12	97.350
7	2	0	102.550
8	2	1	99.650
9	2	3	104.100
10	2	6	101.275
11	2	9	95.850
12	2	12	93.167
13	3	0	104.583
14	3	1	101.200
15	3	3	101.600
16	3	6	100.850
17	3	9	100.925
18	3	12	97.467

The data set represents measurements from three batches. Six replicate assay results, expressed in percent label claim, were measured in months at different storage times, specifically at times 0, 1, 3, 6, 9, and 12 months. Hence, the six time points ($T=6$) were considered at $t = 0, 1, 3, 6, 9,$ and 12 . For all three batches, two replicate response values were missing at times 1, 3, 6, and 9 months. For all three batches, the process begins with a computation of the mean of the replicate assay results across the six replicates. The dataset follows a structure as seen in Table 1.

3.1. Presenting Results from the Analysis

Recall that the confidence limits are of the form $\hat{y}_0 \pm t_{1-\frac{\alpha}{2}} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$. We compute

the 95% confidence interval for the mean of the assay results across the three batches at time 0 month under the fixed-effects model to obtain the values 102.55 and 0.7694

for \hat{y}_0 and $se(\hat{y}_0)$, respectively. Note that $se(\hat{y}_0) = s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$. This yields a

confidence interval of (100.92, 104.18). The 95% confidence intervals for the mean response at the remaining time points are displayed in Table 2. Confidence intervals for the random-effects model are not discussed in this thesis.

A 95% prediction interval that will contain a single future observation from a

population of size n is of the form $\hat{y}_0 \pm t_{1-\frac{\alpha}{2}} s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$. For example if we were

interested in predicting the response from a single future observation at month 0, then a 95% prediction interval is (97.83, 107.26). These intervals are computed for the fixed-effects model and the corresponding prediction intervals for months 1 – 12 are summarized in Table 2. Prediction intervals for the random-effects model are not discussed in this thesis.

Using the cross-sectional structure of the data under the fixed-effects model, the 95%-expectation tolerance of the responses for any batch in the population is of the

form $\bar{Y} \pm K_C S$, where $K_C = \sqrt{1 + \frac{1}{3} t_{0.975(2)}} = 4.303 \sqrt{1 + \frac{1}{3}} = 4.303 \times 1.1547 = 4.970$ and

$t_{\frac{1+\beta}{2}(n-1)}$ is the $100(\frac{1+\beta}{2})$ -percentile in the distribution of $T_{(n-1)}$. At month 0, this interval is (97.78, 108.83). The tolerance intervals for the remaining time points can be found in Table 3.

Recall that under the fixed-effects model, the 95% tolerance interval that contains at least 95% of the population described by a normal distribution is of the form $(\hat{\alpha}_0 + \hat{\alpha}_1 x_0 - g_0 \hat{\sigma}, \hat{\alpha}_0 + \hat{\alpha}_1 x_0 + g_0 \hat{\sigma})$ (Graybill, 1976). Using a longitudinal data structure approach, the 95% tolerance interval that contains at least 95% of the population described by a normal distribution under the fixed-effects model at month 0 is (95.61, 109.49). The tolerance intervals for the remaining time point are given in Table 3.

However, using the longitudinal structure of the data under the random-effects model, the 95%-expectation tolerance interval of the responses for any batch in the population is of the form $\hat{\mu}_i \pm \hat{K}_L \hat{V}$, where the values of \hat{K}_L are given in Table 3. Using the

SAS Software, a regression analysis on the pairs (t_i, Y_{ij}) was done using a random-effects model definition in PROC MIXED. The computation yielded the following values:

$$\begin{aligned}\bar{Y} &= 102.78, & \bar{t} &= 3.5, & W_{tt} &= 110.83, \\ B &= 1.90, & \hat{\alpha} &= 102.55, & \hat{\delta} &= -0.49724 \\ \hat{\sigma}_v &= 0.25771, & \hat{\sigma}_\varepsilon &= 4.16493\end{aligned}$$

At this point \hat{K}_L needs to be determined. Unfortunately, the expression in equation (2.25)

does not readily yield values for \hat{K}_L given R , n , T , C , and A . Instead, the expression is

solved for R by iterating through various values of \hat{K}_L at different time points. Fitting K_L

to polynomial functions of Z using equation (2.26) gives $\hat{R} = 0.67266$ and $\hat{V} = 2.14916$

for equations (2.21) and (2.20) respectively. For example, let's assume that a value of \hat{K}_L

corresponding to $\hat{R} = 0.6727$ is desired to construct a 95% tolerance interval at $t = 0$

months. Let's further assume that $n = 3$, $T = 6$, \bar{t} and W_{tt} are given as aforementioned.

The computer program (see Appendix) gives a list of values of (\hat{K}_L, \hat{R}) from which two

pairs are of our interest: (2.3793, 0.6724) and (2.3821, 0.6673). We conclude from these

two pairs that \hat{K}_L is found in the interval (2.3793, 2.3821) since \hat{R} belongs to the interval

(0.6673, 0.6724). Through a series of iteration of \hat{K}_L from 2.3793 to 2.3821 gives

$(\hat{K}_L, \hat{R}) = (2.3809, 0.6720)$. Therefore, the solution is $\hat{K}_L \approx 2.3809$. The estimated 95%-

expectation tolerance interval for the response Y at month 0 under the random-effects

model is (97.43, 107.66). The variation in the intervals from one time point to another is

due to the change in \hat{K}_L . The tolerance intervals are computed and given in Table 3.

Figure 1 shows a graph of the Jonsson's, Wilks', and Graybill's tolerance intervals.

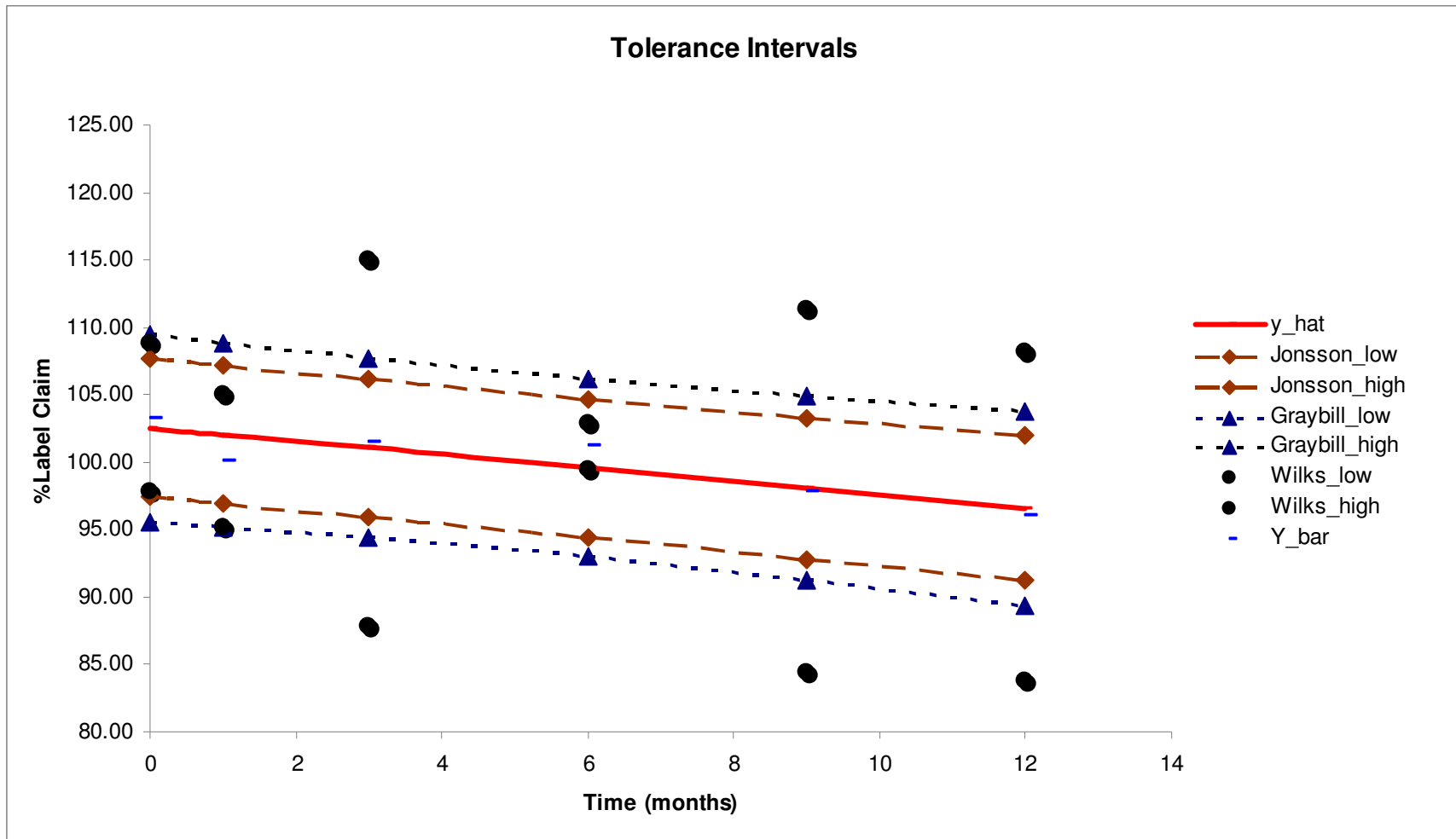
Table 2: 95% Confidence & Prediction Intervals

Month	\bar{Y}	\hat{Y}	Wilks'		95% CI	95% P.I.
			S.D.	S.E. Predicted		
			S	$se(\hat{y}_0)$	Fixed	Fixed
0	103.31	102.55	1.1126	0.7694	(100.92, 104.18)	(97.83, 107.26)
1	100.07	102.05	0.9930	0.6853	(100.60, 103.50)	(97.39, 106.71)
3	101.44	101.06	2.7409	0.5510	(99.89, 102.22)	(96.48, 105.63)
6	101.22	99.56	0.3413	0.5012	(98.50, 100.63)	(95.01, 104.01)
9	97.84	98.07	2.7079	0.6593	(96.67, 99.47)	(93.43, 102.71)
12	96.00	96.58	2.4495	0.9242	(94.62, 98.54)	(91.74, 101.42)

Table 3: 95% Tolerance Intervals

Month	95% T.I.			\hat{K}_L
	Random (Jonsson)	Fixed (Graybill)	Fixed (Wilks)	
0	(97.43, 107.66)	(95.61, 109.49)	(97.78, 108.83)	2.3809
1	(96.94, 107.16)	(95.23, 108.87)	(95.13, 105.00)	2.3768
3	(95.95, 106.16)	(94.43, 107.69)	(87.82, 115.06)	2.3768
6	(94.41, 104.71)	(92.99, 106.13)	(99.52, 102.91)	2.3969
9	(92.83, 103.31)	(91.29, 104.85)	(84.39, 111.30)	2.4381
12	(91.22, 101.94)	(89.39, 103.77)	(83.82, 108.17)	2.465

Figure 1: Tolerance Intervals



3.2. Discussion and Summary of the Analysis

Under the fixed-effects model, the 95% confidence intervals for the mean of the assay results across the three batches are narrower than the 95% prediction intervals. This is mainly justified from the fact that the formula for the 95% prediction intervals include an extra “1” which reflects the additional uncertainty that emanate from a future observation. Tolerance intervals obtained under either the fixed-effects or random-effects model using the longitudinal structure of the data are more appropriate than tolerance intervals computed under the fixed-effects model using the data from a cross-section at each time point. This is an evidence since the former does not use the entire data. Using the longitudinal data structure, the tolerance intervals obtained from Jonsson’s method under the random-effects model are narrower than the tolerance intervals obtained from Graybill’s method under the fixed-effects model. The difference between the two types of tolerance intervals is due to the fact that Graybill’s method for tolerance interval does not incorporate a random effect into the model. We should also notice that there are some fundamental differences between Jonsson’s method, Graybill’s method, and Wilks’ method for tolerance intervals. Jonsson’s method uses the concept of β -expectation tolerance interval under the random-effects model for a longitudinal data structure; the covariance structure is compound symmetry. Graybill’s method uses the concept of β -content tolerance interval under the fixed-effects model for a longitudinal data structure; the covariance structure does assume independence. Wilks’ method uses the concept of β -expectation tolerance interval under the fixed-effects model for a cross-sectional data

structure; no covariance structure is assumed under Wilks' method. The tolerance intervals for Jonsson's, Wilks', and Graybill's methods are shown in Figure 1.

4 Conclusions and Future Research

In this thesis, we described three types of intervals estimates using longitudinal structured data. We first began by introducing the fixed-effects model which treated our population parameter of interest as being fixed with respect to time. We built on the fact that the population parameter of interest was in fact a random-effects parameter as was indicated by most recent methodology. In fact, a major concern in the pharmaceutical industry is that most methodology treats the batches as a fixed effects and therefore ignores the between-batch variability. Inference made based on fixed-effects models are not applicable to future or unobserved batches. Therefore, the use of statistical methods based on a random-effects model (treating intercepts and/or slopes as a random-effects) is more appropriate for establishing the expiry period applicable to future production batches.

We illustrated the confidence intervals for the predicted mean, which are commonly used to establish the expiry period of a drug product. We then described the prediction intervals for a predicted new assay value and the tolerance intervals for a proportion of the population which are sometimes used in risk assessment. However, it is more appropriate to use tolerance intervals in order to assess the risk of product failure at expiry since they provide an interval estimate for the proportion of assay values in the population failing at the expiry period.

We focused our attention to the β -expectation tolerance interval, which requires that $100\beta\%$ of the individual responses from the batches fall between the estimated limits on the average. We chose not to use the other type of tolerance interval, the β -content

tolerance interval, because they are mainly intended for drugs where the risk of adverse effects rapidly increases with an overdose such that even a minor overdose may result in death (Petzold, 2001). Our choice of tolerance interval, the β -expectation tolerance interval, is intended for drugs where the expected outcome of an overdose may only cause discomfort (Petzold, 2001). Since our main interest is on over-the-counter medications, the use of this type of tolerance interval is more appropriate. Furthermore, β -content tolerance intervals appear to be too wide for samples when the within-batch or the between-batch variability is large [Jonsson (2003), Hoffman and Kringle (2005)].

We compared tolerance intervals based on cross-sectional data to tolerance intervals based on longitudinal data. This led us to the conclusion that it is more appropriate to use tolerance intervals derived from longitudinal than those obtained from cross-sectional data since the latter do not use the entire data. This was shown in greater detail by Jonsson (2003).

We finally used the longitudinal data structure to compare tolerance intervals obtained using Jonsson's approach under the random-effects model to tolerance intervals using Graybill's method under the fixed-effects model. We found that the tolerance intervals obtained from Jonsson's method are narrower than the tolerance intervals obtained from Graybill's method. This weakness in Graybill's methodology is due to the fact that it does not incorporate a random effect into the model.

Thus, our findings have led us to a few limitations with the data used for the purpose of the analysis and the methods used for the computation of the random-effects tolerance intervals. In fact, most pharmaceutical stability data are limited when it comes

to the number of batches used. Even though the FDA recommends that at least three batches be used in the analysis of stability data, the minimum number of three batches seems too small for the computation of statistical intervals. The dataset used in our analysis which was performed by Obenchain (1990) falls in the same category. Due to the small number of batches, the intervals computed are wide because of the large within batch variation ($\hat{\sigma}_e = 4.16493$).

Longitudinal-based methods for computing statistical intervals cannot be used without distributional assumptions. However, distribution-free methods are not appropriate for small sample sizes.

The β -expectation tolerance interval presented by Jonsson (2003) is mainly intended for small samples where there is a large within or between subjects variation. This interval was proven to be superior to Wilks' (1941) and Graybill's approaches as the lengths of the intervals are smaller on the average, while the β -expectation property is simultaneously maintained. This is especially the case when the within subjects variation is larger than the between subjects variation.

However, the methodology proposed by Jonsson (2003) needs further improvement. Researchers need to allow for addition of the random slopes in the model since the current approach only allows for random intercepts. This new methodology should also account for the case of unbalanced data. A method could be developed to incorporate the enhanced Jonsson's approach to tolerance intervals in the SAS "MIXED" procedure.

5 List of References

- Johnson, R. A., and Bhattacharyya, G. K. (2006). *Statistics: Principles and Methods*. (Fifth Edition). John Wiley & Sons, Inc., Hoboken, NJ.
- Jonsson, R. (2003). "A Longitudinal Approach for Constructing β -Expectation Tolerance Intervals", *Journal of Biopharmaceutical Statistics*, 13, 2, 307-325.
- Mongomery, D.C., and Runger, G.C. (2003). *Applied Statistics and Probability For Engineers*. (Third Edition). John Wiley & Sons, Inc., New York, NY.
- Murphy, J.R., and Hofer, J.D. (2002). "Establishing Shelf Life, Expiry Limits, and Release Limits", *Drug Information Journal*, 36, 769-781.
- Littell, R.C., Milliken, G.A., Stroup, W.W., and Wolfinger, R.D. (1996). *SAS[®] System for Mixed Models*, SAS Institute Inc., Cary, NC.
- Chow, S., and Liu, J. (1995). *Statistical Design and Analysis in Pharmaceutical Science*. Marcel Dekker, New York, NY.
- Hahn, G.J., and Meeker, W.Q. (1991). *Statistical Intervals: A Guide for Practitioners*. John Wiley & Sons, Inc., New York, NY.
- Rode, R.A. (1986). *The Use of Box-Cox Transformations in the Development of Multivariate Tolerance Regions with Applications to Clinical Chemistry*. Virginia Commonwealth University Press, Richmond, VA.
- Kempthorne, O., and Folks, L. (1971). *Probability, Statistics, And Data Analysis*. The Iowa State University Press, Ames, IO.
- Nickens, D. J. (1998). Using tolerance limits to evaluate laboratory data. *Drug Inf. J.*

32:261–269.

Sommerville, P. N. (1958). Tables for obtaining non-parametric tolerance limits. *Ann.*

Math. Stat. 29:599–601.

Graybill, F. A. (1976). *Theory and Applications of the Linear Model*. Wadsworth

Publishing Company, Inc., Belmont, CA.

6 Appendices

6.1. SAS Code

```

/*****
*****
*****
*****
Kakotan Sanogo
Under the Guidance of Drs. Jessica M. Ketchum,
Charles W. Kish,
and Ramakrishnan Viswanathan
Department of Biostatistics, Virginia Commonwealth
University
Division of Statistics, Wyeth Consumer Healthcare
Modified on: 11-Dec-2008
Dataset: Obenchain (1990)
References: Analysis of Messy Data, Miliken & Johnson, 2002
The Theory and Application of the Linear Model,
Graybill, 1976
A Longitudinal Approach for Constructing beta-
Expectation
Tolerance Intervals, Jonsson, 2003
Determination of Sample Sizes for Setting Tolerance
Limits, Wilks, 1941
*****
*****
*****
*****/

data rc;
  input batch month r1-r6;
  array r{6};
  monthc = month;
  drop i r1-r6;
  do i = 1 to 6;
    y = r{i};
    if (y ^= .) then output;
  end;
  datalines;
1 0 101.2 103.3 103.3 102.1 104.4 102.4
1 1 98.8 99.4 99.7 99.5 . .
1 3 98.4 99.0 97.3 99.8 . .
1 6 101.5 100.2 101.7 102.7 . .
1 9 96.3 97.2 97.2 96.3 . .
1 12 97.3 97.9 96.8 97.7 97.7 96.7
2 0 102.6 102.7 102.4 102.1 102.9 102.6
2 1 99.1 99.0 99.9 100.6 . .

```

```

2 3 105.7 103.3 103.4 104.0 . .
2 6 101.3 101.5 100.9 101.4 . .
2 9 94.1 96.5 97.2 95.6 . .
2 12 93.1 92.8 95.4 92.5 92.2 93.0
3 0 105.1 103.9 106.1 104.1 103.7 104.6
3 1 102.2 102.0 100.8 99.8 . .
3 3 101.2 101.8 100.8 102.6 . .
3 6 101.1 102.0 100.1 100.2 . .
3 9 100.9 99.5 102.5 100.8 . .
3 12 97.8 98.3 96.9 98.4 96.9 96.5
;
proc sort; by batch month;
proc means; by batch month; var y; output out = ymeans mean = y;
proc print data = ymeans;
run;
data test;
set ymeans(keep=batch month y);
proc print;
run;

*****
*****
*****
COMPUTATION OF CONFIDENCE INTERVALS: FIXED-EFFECTS MODEL
*****
*****
***** Method ONE;
proc mixed data = test;
class batch;
model y = month /solution outp = out1 outpm = out2 cl;
title '95% Confidence Intervals for the Fixed-Effects Model';
run;
proc print data = out2;
run;

***** Method TWO;
proc mixed data = test;
class batch;
model y = month /solution outp = out1 outpm = out2 cl;
title '95% Confidence Intervals for the Fixed-Effects Model';
estimate "time0" intercept 1 month 0/ cl df=16;
estimate "time1" intercept 1 month 1/ cl df=16;
estimate "time3" intercept 1 month 3/ cl df=16;
estimate "time6" intercept 1 month 6/ cl df=16;
estimate "time9" intercept 1 month 9/ cl df=16;
estimate "time12" intercept 1 month 12/ cl df=16;
run;

**** Method THREE;
proc glm data=test ;
model y=month / SS3 clm ; *clm Prints 95% confidence

```

```

                                intervals for the
mean of each observation ;
    title '95% Prediction Intervals for the Fixed-Effects Model';
    estimate "time0" intercept 1 month 0 ;
    estimate "time1" intercept 1 month 1 ;
    estimate "time3" intercept 1 month 3;
    estimate "time6" intercept 1 month 6;
    estimate "time9" intercept 1 month 9;
    estimate "time12" intercept 1 month 12;
run;

*****
*****
*****
COMPUTATION OF CONFIDENCE INTERVALS: RANDOM-EFFECTS MODEL
*****
*****
*****Method ONE;
proc mixed data = test;
    class batch;
    model y = month / s cl outpm=predm ;
    random int month/ sub=batch s g v;
    title '95% Confidence Intervals for the Random-Effects Model';
run;
proc print data = predm;
run;

*****Method TWO;
proc mixed data = test;
class batch;
    model y = month /solution outp = out1 outpm = out2 cl ;
    random intercept month / sub= batch type = un solution v vcorr;
    title '95% Confidence Intervals for the Random-Effects Model';
    estimate "time0" intercept 1 month 0/ cl df=12;
    estimate "time1" intercept 1 month 1/ cl df=12;
    estimate "time3" intercept 1 month 3/ cl df=12;
    estimate "time6" intercept 1 month 6/ cl df=12;
    estimate "time9" intercept 1 month 9/ cl df=12;
    estimate "time12" intercept 1 month 12/ cl df=12;
run;

/*****
*****
COMPUTATION OF PREDICTION INTERVALS: FIXED-EFFECTS MODEL
*****
*****/
proc glm data=test;
    model y=month / SS3 cli ; *cli Prints 95% prediction
                                intervals for the mean
of each observation ;
title '95% Prediction Intervals for the Fixed-Effects Model';

```

```

estimate "time0" intercept 1 month 0 ;
estimate "time1" intercept 1 month 1 ;
estimate "time3" intercept 1 month 3;
estimate "time6" intercept 1 month 6;
estimate "time9" intercept 1 month 9;
estimate "time12" intercept 1 month 12;

run;

/*****
*****
COMPUTATION OF PREDICTION INTERVALS: RANDOM-EFFECTS
*****
*****/

%let _time=Month;
%let _alpha=0.05;

proc mixed data =test ;

class batch ;
model y = &_amp;_time / s cl outpm=predm ;
      random int &_amp;_time/ sub=batch s g v;
      ods output solutionf=solutionf;
      ods output solutionr=solutionr;
      ods output covparms=covparms;
      ods output dimensions=dim;
      ods output tests3=tests3;
      ods output nobs=nobs;
      ods output ClassLevels=class;

run ;

/* PRINT DATASETS CREATED FROM PROC MIXED */
data dim2;
  set dim;
  do i=1 to 18 by 1;
  if descr='Subjects' then
  do;
    n_subj=value;
    df_intercept = n_subj - 1;
    df_slope = n_subj - 1;
    call symput('&_n_subj', n_subj);
    output;
  end;
  keep n_subj df_intercept df_slope;
end;
%put &_n_subj;

run;
proc print data=dim2;run;

*****
CREATE A DATASET WITH NUMBER OBS USED IN ANALYSIS AND
DEGREES OF FREEDOM FOR MSE

```



```

*****
**;
  data nobs2;
    set nobs;
    do i=1 to 18 by 1;
      if label="Number of Observations Used" then
        do;
          n_tot=n;
          df_err = n_tot - (2*&n_subj);
          output;
          end;
        keep n_tot df_err ;
        end;
      run;
proc print data=nobs2;run;

*****
*
  CREATE DATASET OF INTERCEPT ESTIMATES USING PREDM DATASET
  & MERGE WITH PREDM DATASET

*****
;
  data intercept;
    set predm;
    if &_time=0 then
      do;
        intercept=pred;
        output;
        end;
      keep batch &_time intercept;
    run;
proc print data=intercept;run;

proc sort data=intercept out=intercept; by batch;
proc sort data=predm out=predm; by batch;

  data predm_int;
    merge predm intercept;
    by batch;
  run;
proc print data=predm_int;run;

  data predm_int2;
    set predm_int(rename=(intercept=temp_int));
    retain x .;
    intercept=temp_int;
    if temp_int ne ' ' then x = intercept;
    else if temp_int = ' ' then intercept = x;
    drop x temp_int;
  run;
proc print data=predm_int2;run;

```

```

*****
*
  CREATE DATASET OF SLOPE ESTIMATES USING PREDM DATASET
  & MERGE WITH PREDM / INTERCEPT DATASET
*****
;
  data slope;
    set predm_int2;
    if &_amp;_time ne 0 then
      do;
        slope = (pred - intercept) / &_amp;_time;
        output;
      end;
    keep batch &_amp;_time slope;
  run;
proc print data=slope;run;

  proc sort data=slope;
    by batch descending &_amp;_time;
  run;

  proc sort data=predm_int2;
    by batch descending &_amp;_time;
  run;

  data predm_int_slope;
    merge predm_int2 slope;
    by batch;
  run;
proc print data=predm_int_slope;run;

  data predm_int_slope2;
    set predm_int_slope(rename=(slope=temp_slope));
    retain x .;
    slope=temp_slope;
    if temp_slope ne ' ' then x = slope;
    else if temp_slope = ' ' then slope = x;
    drop x temp_slope;
  run;
proc print data=predm_int_slope2;run;

*****
*
  CREATE A DATASET WITH MSE & MERGE WITH PREDICTED VALUES DATASET
*****
;
  proc print data=covparms;
  run;

```

```

data covparms2;
  set covparms;
  retain var_intercept var_slope var_err;
  do i=1 to 18 by 1;
    if covparm='Intercept' then var_intercept=estimate;
    else if covparm="&_time" then var_slope=estimate;
    else if covparm='Residual' then
      do;
        var_err=estimate;
        output;
      end;
    drop covparm estimate;
  end;
run;
proc print data=covparms2;run;

*****
*
*      COMBINE PREDM/INTERCEPT/SLOPE, COVPARM, & NOBS ESTIMATES
*****
;
  proc sort data=predm_int_slope;
  by batch descending &_time;
run;

data reg_est;
  merge predm_int_slope2 covparms2 dim2 nob2;
  *by batch;
run;
proc print data=reg_est;run;

*****
*****
  RESTRUCTURE RANDOM EFFECT DEVIATES DATASET
*****
;
  proc sort data=solutionr out= solutionr;by batch ;

proc print data=solutionr;
run;

data solutionr2;
  set solutionr;
  by batch;
  retain effect_i effect_s .;
  if effect="Intercept" then effect_i=estimate;
  else if effect ne "Intercept" then effect_s=estimate;
  if last.batch then output;

```

```

        keep batch effect_i effect_s;
run;
proc print data=solutionr2;run;

*****
*
      COMBINE PREDM/INTERCEPT/SLOPE/COVPARAM/NOBS WITH RANDOM EFFECT
DEVIATES ESTIMATES

*****
;
proc sort data=solutionr2 out= solutionr2;by batch;
proc sort data=reg_est out= reg_est;by batch;
data reg_est2;
  merge reg_est solutionr2;
  by batch;
run;
proc print data=reg_est2;run;

*****
*
      COMPUTE LOWER PREDITION LIMIT

*****
;
data reg_est3;
  set reg_est2;

  *** standard error of the predicted value based on the mixed
model;
  *** assumes no covariance between intercept & slope;
  var_pred_vc=var_intercept + ((&_time**2)*var_slope) + var_err;

  *** degrees of freedom based on satterthwaites approximation;
  num= ((var_intercept + ((&_time**2)*var_slope) + var_err)**2);
  denom= ((var_intercept**2)/df_intercept) +
  (((&_time**2)*var_slope)**2) / df_slope) + ((var_err**2) / df_err);

  df_sw = num/denom;

  fcrit=finv(1-&_alpha,2,df_sw);

  * f-distribution check;
  fcrit_ck=finv(1-&_alpha,2,70);
  scheffe_pct_pt = sqrt(2*fcrit_ck);

  piw=sqrt(2*fcrit)*sqrt(var_pred_vc);

  *** compute 2-sided upper pred. limits about a future obs.;
  pred_chk=intercept + (slope*&_time);

  lpl = pred - piw;

```

```

        upl = pred + piw;
        int_width=upl-lpl;
run;

proc print data=reg_est3;run;

*****
*
  CREATE LISTING OF REGRESSION ESTIMATES AND RELEASE LIMITS (LOWER &
UPPER)
*****
;

ods select all;

proc print data=reg_est3 label;
      var &_time batch intercept SLOPE pred var_intercept
var_slope var_err
      df_intercept df_slope df_err fcrit lpl upl int_width df_sw
var_pred_vc piw;
run;

/*****
***
**
  COMPUTATION OF TOLERANCE INTERVALS: Wilks' Method
*****
**
*****/
data _stats;
input y1 y2 y3@@;
datalines;
102.783      102.55      104.583
99.35  99.65      101.2
98.625      104.1      101.6
101.525      101.275  100.85
96.75  95.85      100.925
97.35  93.167      97.467
;
run;
proc print data=_stats;run;

data wilks; set _stats;

*Computing the mean batch;
y_bar = mean(y1,y2,y3);

* Computing the standard deviation;

```

```

SD = std(y1,y2,y3);

K = sqrt(1+1/3)*(4.303); * where n=3 batches and 4.303 is the
97.5 percentile
                                in a t distribution
with 2 degrees of freedom;
    LTL = y_bar - K*SD;
    UTL = y_bar + K*SD;
run;

proc print data= wilks;
    var LTL UTL SD y_bar;
run;

/*****
*****
COMPUTATION OF TOLERANCE INTERVALS: Graybill's Method
*****
*****/
proc sort data=rc; by batch month;
proc means; by batch month; var y; output out = ymeans mean = y;
proc print data = ymeans;
run;
data test;
    set ymeans(keep=batch month y);
    proc print;
run;
proc glm data=test;
    model y=month / ss3;
    output out=predcheck predicted=pred h=h;
run;

proc print data=predcheck;
run;

data work;
    length method $9;
    input method $& alpha p pred mse df h expiry;
    /* the value of "pred", "mse", "df", and "h" are
       obtained from the "glm" procedure output and "pred" and "h"
       vary for different values of "expiry" (from 0 to 12 month).
       P=0.05 represent the 95% of tolerance point with (1-alpha)=95%
confidence.
       The dataset below computes the Graybill's Tolerance Interval
       at month 0.
    */
    cards;
    Graybill .05 .05 102.55 4.3582101 16 0.13584 0
    ;
run;

proc print data=work;
run;
data stats;set work;

```

```

rootmse=mse**.5;
A=h**.5;

    np=probit(1-(p/2));
    deltau=-np/A;
    tcritu=tinv(alpha/2, df, deltau);
    gpu= - A * tcritu;
    UTW=gpu * rootmse;
    UTL = pred + (gpu * rootmse);

    deltal=np/A;
    tcritl=tinv(1-(alpha/2), df, deltal);
    gpl= A * tcritl;
    LTW=gpl * rootmse;
    LTL = pred - (gpl * rootmse);
run;

proc print data=stats;
    var LTL UTL;
run;

/*****
*****
COMPUTATION OF TOLERANCE INTERVALS: Jonsson's Method
*****
*****/

/*****
***
**
START GLOBAL MACRO
*****
**
*****/
* First run the global macro and then change the value of _time
to
get the tolerance intervals;
%macro bigone(_time, _n, _TotTimes);

%do _t=0 %to &_time %by 1;
    data a;
    set test;
    i = month;
    i2 = i*i;
    Y2 = Y*Y;
    iY = i*Y;
    run;
proc sort;
    by batch;
run;
proc means noprint sum;

```

```

        var i Y i2 Y2 iY;
        by batch;
        output out = sas1 sum = si sY si2 sY2 siY;
run;
proc print data = sas1;
run;
data b;
    set sas1;
    mt = si/&_TotTimes;      mYj = sY/&_TotTimes;      mYj2 =
mYj*mYj;      Wttj = si2-si*si/&_TotTimes;
    WtYj = siY-si*sY/&_TotTimes;      WYjYj = sY2-sY*sY/&_TotTimes;
    bj = WtYj/Wttj;      aj = mYj-bj*mt;
run;
proc means noprint sum;
    var bj aj mYj mYj2 Wttj WYjYj WtYj;
    output out = sas2 sum = sbj saj smYj smYj2 sWttj sWYjYj
sWtYj;
run;
proc print data= sas2;
run;

/* totals */
data c;
    set sas2;
    b = sbj/&_n; a = saj/&_n; Wtt = sWttj/&_n; BYY = smYj2-
smYj*smYj/&_n;
    vare = (sWYjYj-b*sWtYj)/(&_n*(&_TotTimes-1)-1); varu =
BYY/(&_n-1)-vare/&_TotTimes;
    R = vare*(1-1/&_TotTimes)*(&_n*(&_TotTimes-1)-
3)/(vare+varu)/(&_n*(&_TotTimes-1)-1);
    Z = (1-R)*(1-R)/(&_n-1)+R*R/(&_n*(&_TotTimes-1)-1);
    V = sqrt(vare+varu)/(1-Z/4);
    call symput('R', R);run;
proc print ;
    var a b Z V R vare varu Wtt BYY;
run;

/*****
*****
*****
*****
THIS MACRO CREATES K_hat for a given TIME POINT
*****
*****
*****
*****/

data d;
meanT = 2.0;      Wtt = 110.83;      beta = 0.95;

const = gamma(1/2)*sqrt(2);
do K = 2.0 to 3.0 by 0.0001;
    C = (&_time - meanT)**2/&_n/(1-1/&_TotTimes)/Wtt;

```



```

        AA = 2*(const*(exp(K*K/2))*(2*probnorm(K)-1-
beta)/K-1/&n)/K/K;
        R_hat = ((&n*(&TotTimes-1)-1)*(1-(&n-1)*C/K/K)-
sqrt((&n*(&TotTimes-1)-1)**2*(1-(&n-1)*C/K/K)**2-(&n*&TotTimes-
2)*(&n*(&TotTimes-1)-1)*(1-(&n-1)*AA)))/(&n*&TotTimes-2);
        dif_R=abs(R_hat-&R);
        output;
    end;
run;

data e; set d; if dif_R=. then delete; run;

proc sort data=e; by dif_R; run;

data f; set e;
    if &n ne 1 then delete;
    call symput("K", K);run;
run;

proc print data = f;
    var K;
    title ' K_hat at Time &time when beta = 0.95';
run;

/*****
*****
*****
***** THIS MACRO CREATES TOLERANCE INTERVALS FOR EACH TIME POINT
*****
*****
*****/
data ti;
    set c d;

    low_ti = (a+b*&_time)-&K*V;
    up_ti = (a+b*&_time)+&K*V;
proc sort; by low_ti up_ti;
run;

data tol_int; set ti;
    if low_ti=. then delete;
    else if up_ti=. then delete;
run;

proc print data=tol_int;
    var low_ti up_ti;
run;

run;quit;
%end;
%mend;

```

```

%bigone(0,3,6); * change the value of _time in bigone(_time, _n,
_TotTimes) to
                obtain the values of Kl and the tolerance intervals;

/*****
*****
END GLOBAL MACRO
*****
*****/

/*****
*****
PLOTS OF TOLERANCE INTERVALS
*****
*****/

data TI;
input t y Jonsson_low Jonsson_high Graybill_low Graybill_high Wilks_low
Wilks_high y_bar;
cards;
0      102.55      97.43 107.66      95.61 109.49      97.7771
      108.834      103.31
1      102.05      96.94 107.16      95.23 108.87      95.1333      105
100.07
3      101.06      95.95 106.16      94.43 107.69      87.8229
      115.06      101.44
6      99.56 94.41 104.71      92.99 106.13      99.5211      102.912
      101.22
9      98.07 92.83 103.31      91.29 104.85      84.387      111.296
      97.84
12     96.58 91.22 101.94      89.39 103.77      83.8237      108.166
      96
;
* The following exports the data into microsoft excel;
ods html file = 'C:\Documents and
Settings\sanogok\Desktop\THESIS\thesis.xls' ;
proc print data =TI;
title 'Tolerance Intervals' ;
run ;
ods html close ;
* The excel statemtentns for creating the plot are the following ;

* For y_hat
=SERIES(Sheet1!$B$1, Sheet1!$A$2:$A$7, Sheet1!$B$2:$B$7,1);

* For Jonsson_low
=SERIES(Sheet1!$C$1, Sheet1!$A$2:$A$7, Sheet1!$C$2:$C$7,1);

* For Jonsson_high
=SERIES(Sheet1!$D$1, Sheet1!$A$2:$A$7, Sheet1!$D$2:$D$7,1);

* For Graybill_low
=SERIES(Sheet1!$E$1, Sheet1!$A$2:$A$7, Sheet1!$E$2:$E$7,1);

```

```
* For Graybill_high
=SERIES(Sheet1!$F$1, Sheet1!$A$2:$A$7, Sheet1!$F$2:$F$7, 1);

* For y_bar
=SERIES(Sheet1!$I$1, Sheet1!$A$2:$A$7, Sheet1!$I$2:$I$7, 1);

* For Wilks_low
=SERIES(Sheet1!$G$1, Sheet1!$A$2:$A$7, Sheet1!$G$2:$G$7, 5);

* For Wilks_high
=SERIES(Sheet1!$H$1, Sheet1!$A$2:$A$7, Sheet1!$H$2:$H$7, 5);

* t represents the x-axis and y the y-axis;
```

VITA

Kakotan Sanogo was born on October 8, 1980 in Koulikoro, Mali (West Africa). He graduated from Lycee Prosper Kamara (Bamako, Mali) in 1998 and from Lycee d'Etat de l'Estuaire (Libreville, Gabon – Central Africa) in 1999. He received his Bachelors of Science in Computer Engineering from Minnesota State University – Mankato in May of 2004. At MNSU he was recognized for his academic and extracurricular activities achievements receiving the Who's Who Among the Students in American Universities and Colleges Award, as well as scholarships from the departments of Mathematics and Electrical and Computer Engineering. In August of 2004, he moved to Ruston, Louisiana, to begin his graduate studies in Biomedical Engineering at Louisiana Tech University. While in Minnesota and Louisiana, he spent two and a half years working at local nursing homes as a Certified Nursing Assistant in an attempt to serve the communities in which he lived. On his way to Louisiana, he learned about the fascinating field of Biostatistics upon a short stop over his cousin's, Dr. Oumar Sy, a Senior Research Biostatistician who was at that time preparing his dissertation defense for the Doctor of Philosophy degree in Mathematical Statistics at Kansas State University. He immediately applied to the graduate program in Biostatistics at Virginia Commonwealth University which he started in August of 2005. He spent the first year of his biostatistics graduate career as a teaching assistant to biostatistics classes taught to medical and graduate students, followed by a year of consulting with the Virginia Commonwealth University Technology Services; he subsequently spent a year-long internship at Wyeth

Consumer Healthcare where he learned statistical methods used in the drug development process and about life in corporate America. This was followed by a two-month adjunct faculty position in the department of Pediatric Dentistry where he assisted a faculty member with statistical programming for her grant. He currently is employed as a Clinical Data Analyst for the division of Epidemiology and Infectious Diseases within the department of Internal Medicine at the Virginia Commonwealth University Health System. He practices Taekwondo, a form of martial art, and reads business and self-improvement books in his spare time. He will get married on December 27th 2008.