

2018

Comparison of the Performance of Simple Linear Regression and Quantile Regression with Non-Normal Data: A Simulation Study

Marjorie Howard
University of South Carolina

Follow this and additional works at: <https://scholarcommons.sc.edu/etd>

 Part of the [Biostatistics Commons](#)

Recommended Citation

Howard, M.(2018). *Comparison of the Performance of Simple Linear Regression and Quantile Regression with Non-Normal Data: A Simulation Study*. (Doctoral dissertation). Retrieved from <https://scholarcommons.sc.edu/etd/4517>

This Open Access Dissertation is brought to you by Scholar Commons. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of Scholar Commons. For more information, please contact dillarda@mailbox.sc.edu.

Comparison of the Performance of Simple Linear Regression and Quantile Regression
with Non-Normal Data: A Simulation Study

by

Marjorie Howard

Bachelor of Arts
Mercer University, 2010

Submitted in Partial Fulfillment of the Requirements
For the Degree of Master of Science in Public Health in
Biostatistics

The Norman J. Arnold School of Public Health

University of South Carolina

2018

Accepted by:

Andrew Ortaglia, Director of Thesis

James Hussey, Reader

Michael Wirth, Reader

Cheryl L. Addy, Vice Provost and Dean of the Graduate School

© Copyright by Marjorie Howard, 2018

All Rights Reserved.

ABSTRACT

Linear regression is a widely used method for analysis that is well understood across a wide variety of disciplines. In order to use linear regression, a number of assumptions must be met. These assumptions, specifically normality and homoscedasticity of the error distribution can at best be met only approximately with real data. Quantile regression requires fewer assumptions, which offers a potential advantage over linear regression. In this simulation study, we compare the performance of linear (least squares) regression to quantile regression when these assumptions are violated, in order to investigate under what conditions quantile regression becomes the more advantageous method of analysis. Statistical power and coverage percentage were calculated for all simulations, and potential bias was investigated for both quantile regression and linear regression. When errors are skewed, there is a threshold at which quantile regression surpasses linear regression in statistical power. When heteroscedasticity is introduced, linear regression does not accurately describe the relationship between predictor and response variables at the tails of the conditional distribution. When errors are both skewed and heteroscedastic, quantile regression performs drastically better than linear regression. Coverage percentage in linear regression not only suffers in this case, but also linear regression yields misleading results.

TABLE OF CONTENTS

Abstract.....	iii
List of Tables	v
List of Figures.....	viii
Chapter 1. Introduction.....	1
Chapter 2. Background	2
2.1. Linear Regression	2
2.2. Quantile Regression.....	4
Chapter 3. Simulation Design.....	7
3.1. Scenario 1.....	7
3.2. Scenario 2.....	10
3.3. Scenario 3.....	11
Chapter 4. Results.....	35
4.1. Scenario 1.....	35
4.2. Scenario 2.....	37
4.3. Scenario 3.....	39
Chapter 5. Conclusions.....	50
References.....	54

LIST OF TABLES

Table 3.1 Simulation Scenario 1: Error Distributions by Skewness.....	15
Table 3.2 Estimates for Mean(SLR) and Median(QR) when Linear Regression Assumptions Hold.....	16
Table 3.3 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 0.5.....	16
Table 3.4 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 1.....	17
Table 3.5 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 2.....	17
Table 3.6 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 2.5.....	18
Table 3.7 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 3.....	18
Table 3.8 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 4.....	19
Table 3.9 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 6.....	19
Table 3.10 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -0.5	20
Table 3.11 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -1	20
Table 3.12 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -2	21
Table 3.13 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -2.5	21

Table 3.14 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -3	22
Table 3.15 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -4	22
Table 3.16 Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -6	23
Table 3.17 Type I Error Rate for SLR and Median QR with Slope Coefficient of 0	24
Table 3.18 Simulation Scenario 2 Error Distributions.....	24
Table 3.19 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5.....	25
Table 3.20 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2.....	25
Table 3.21 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3.....	26
Table 3.22 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5.....	26
Table 3.23 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2.....	27
Table 3.24 Comparison of SLR(Mean) and QR(<i>p</i> th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3.....	27

Table 3.25 Comparison of SLR(Mean) and QR(p th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5.....	28
Table 3.26 Comparison of SLR(Mean) and QR(p th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2.....	28
Table 3.27 Comparison of SLR(Mean) and QR(p th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3.....	29
Table 3.28 Estimates for SLR and QR(p th percentile) with Skewed Errors and Heteroscedasticity	29

LIST OF FIGURES

Figure 3.1 Scenario 1: Histogram of Error with Skewness 0.5 Compared with Normal Curve.....	30
Figure 3.2 Scenario 1: Histogram of Error with Skewness 1 Compared with Normal Curve.....	30
Figure 3.3 Scenario 1: Histogram of Error with Skewness 2 Compared with Normal Curve.....	31
Figure 3.4 Scenario 1: Histogram of Error with Skewness 2.5 Compared with Normal Curve.....	31
Figure 3.5 Scenario 1: Histogram of Error with Skewness 3 Compared with Normal Curve.....	32
Figure 3.6 Scenario 1: Histogram of Error with Skewness 4 Compared with Normal Curve.....	32
Figure 3.7 Scenario 1: Histogram of Error with Skewness 6 Compared with Normal Curve.....	33
Figure 3.8 Scenario 2: Histogram of Error with Skewness and Heteroscedasticity at ($x=-2$)	33
Figure 3.9 Scenario 2: Histogram of Error with Skewness and Heteroscedasticity at ($x=0$).....	34
Figure 3.10 Scenario 2: Histogram of Error with Skewness and Heteroscedasticity at ($x=2$).....	34
Figure 4.1 Coverage Percentage by Skewness and Sample Size for SLR (Slope=0.4).....	41
Figure 4.2 Coverage Percentage by Skewness and Sample Size for Median QR (Slope=0.4)	42
Figure 4.3 Power Comparison between SLR and Median QR (Slope=0.4) by Sample Size and Skewness of Error	43

Figure 4.4 Comparison of Mean Absolute Bias(%) between SLR and Median QR by Sample Size and Skewness of Error	44
Figure 4.5 Type I Error Rate of SLR and Median QR by Sample Size and Skewness of Error	45
Figure 4.6 SLR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity	46
Figure 4.7 10 th Percentile QR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity.....	47
Figure 4.8 90 th Percentile QR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity.....	48
Figure 4.9 Type I Error Rate of SLR by Sample Size and Level of Heteroscedasticity ...	49

CHAPTER 1

INTRODUCTION

Linear regression is a common and well understood method for modelling and predicting the mean of a continuous response variable, conditional on a set of predictor variables. However, inference from linear regression requires specific assumptions be made about the distribution of the error. Quantile regression has become an increasingly popular alternative to linear regression as it can be used to predict any desired percentile of a continuous response variable, conditional on a set of predictor variables, without making assumptions about the distribution of the error. In this paper, we present a simulation experiment to compare the performance of simple linear regression (SLR) to simple quantile regression (QR) when the conditional normality and equal variance assumptions required for linear regression do not hold. Our goal is to provide recommendations for when quantile regression might be preferred over linear regression and vice versa.

CHAPTER 2

BACKGROUND

2.1 LINEAR REGRESSION

Linear regression is a common statistical method which is used to gain understanding of the relationship between a continuous response variable and a set of predictor variables. This model estimates the mean of the response variable, conditional on fixed values for the predictor variables. The mean of the outcome, conditional on these predictors, is often referred to as the conditional mean. In the case of simple linear regression, one predictor variable (generally denoted by x) is used to model the response variable (denoted by y). For $i = 1, 2, \dots, n$ independent observations of a given independent variable x_i with an outcome variable y_i , the model is defined as $y_i|x_i = \beta_0 + \beta_1 * x_i + \varepsilon_i$, where β_0 and β_1 are regression parameters, and indicate the mean response when $x_i = 0$, and the change in the mean response associated with a one unit increase in the independent variable x_i , respectively (Kleinbaum, 2008). The error term is represented by ε_i and is defined as the difference from the population mean response, associated with the corresponding x_i value, for an individual observation.

In linear regression, the regression parameters can be estimated using the ordinary least squares method, which estimates regression parameters based upon minimization of the sum of squared residuals. Residuals are defined as the difference between the observed response and the predicted response. The ordinary least squares method yields

the following parameter estimates for β_0 and β_1 , using \bar{x} to represent the mean value of the predictor variable, \bar{y} to represent the mean value of the response variable, and hats above the parameters, $\widehat{\beta}_0$ and $\widehat{\beta}_1$ to denote the predicted or estimated values for these parameters (Kleinbaum, 2008):

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$$

Inference from linear regression relies upon a number of assumptions. As with many statistical models, linear regression models assume that each of the observations are independent of one another. It is also assumed that the mean of the response variable can be represented by a linear combination of predictor variables. In addition, linear regression models assume that the distribution of the random error is normal with mean zero and that the variance of the response variable is the same for any value of the predictor variable (e.g. $\varepsilon \sim N(0, \sigma^2)$), commonly referred to as homogenous variance. This assumption forms the basis of standard error calculations, as the variance of the conditional distribution is estimated as a single parameter estimate, $\frac{\sum_{i=1}^n (y_i - \hat{y}_i)}{n-2}$ (Pagano, 2000). If these assumptions hold true, the distribution of $y_i | x_i$ can be specified as $N(\beta_0 + \beta_1 * x_i, \sigma^2)$. With the assumption of normality, ordinary least squares methods yield the same parameter estimates as maximum likelihood estimation.

If all of the assumptions are met, linear regression provides a widely accepted and easily understood description of the association between two variables. That is, that any

one-unit increase in the predictor variable is associated with a corresponding change of β_1 in the mean of the response variable.

However, none of these assumptions can ever be met exactly with real data. When these assumptions are violated, the results and inference from linear regression can be impacted. For example, when errors are heteroscedastic, the variance in the response variable may be over or underestimated at times, which impacts standard error calculations, and thus confidence intervals and statistical power. In addition, with heteroscedasticity in errors, the parameters estimated with linear regression might be unable to accurately describe the full range of the associations across different percentiles of the conditional distribution of the outcome. The point estimates at the mean may over or underestimate the true association at the tails of the conditional distribution of the response variable.

2.2 QUANTILE REGRESSION

Quantile regression is similar to linear regression in that it is used to gain an understanding of how a set of predictor variables are related to a continuous response variable; however, there are several differences in the calculation, assumptions, and interpretation when comparing quantile regression to linear regression. Koenker and Basset (1978), introduced quantile regression as a way to model the quantile of a response variable, e.g. the median, conditional on a certain value of a set of predictors. As the quantile is estimated conditional on the set of predictors, it commonly referred to as the conditional quantile. In the field of public health, researchers are often most interested in the more extreme values of the conditional distribution rather than the mean, making quantile regression an excellent method of analysis in this field. For example, in studies

of factors which may impact birthweight, researchers are primarily interested in the lower tails of the distribution of birthweight, as low birthweight can have medical and financial long-term effects. Abrevaya (2001) used quantile regression to suggest that the effect certain factors have upon a child's birthweight can vary depending on the quantile of birthweight.

For a sample of $i = 1, \dots, n$ observations of a predictor variable x_i and a response variable y_i with conditional probability distribution $F_Y(y)$ given x , the conditional quantile p is defined as $Q_{y|x}(p) = \text{infimum}\{y : F_Y(y) \geq p\}$ (Hao & Naiman, 2007). Rather than minimizing the sum of the squared residual as in linear regression, estimation in quantile regression is done by minimizing a weighted sum of absolute residuals. The equation below shows the quantile regression model as specified for the p th quantile,

$$Q_{y_i|x_i}(p) = \beta_0^{(p)} + \beta_1^{(p)} * x_i + \varepsilon_i^{(p)}$$

Here, $\beta_0^{(p)}$ and $\beta_1^{(p)}$ represent the p th quantile of the response variable when $x = 0$ and the change in the p^{th} quantile for a one unit increase in the predictor variable, respectively. If we assume this linear relationship between the predictor and response variables, $\beta_0^{(p)}$ and $\beta_1^{(p)}$ can be calculated as solutions to the minimization of the weighted sum of absolute residuals, as given by (Koenker & Bassett, 1978)

$$\min_{\beta_0, \beta_1} \left(p * \sum_{y_i \geq \widehat{\beta}_0^{(p)} + \widehat{\beta}_1^{(p)} * x_i} |y_i - \widehat{\beta}_0^{(p)} - \widehat{\beta}_1^{(p)} x_i| + (1 - p) * \sum_{y_i < \widehat{\beta}_0^{(p)} + \widehat{\beta}_1^{(p)} * x_i} |y_i - \widehat{\beta}_0^{(p)} - \widehat{\beta}_1^{(p)} x_i| \right)$$

In this process, fitted values that underpredict an observed value are given a weight of $1 - p$, and those that overpredict are given a weight of p . Through this process of applying weights to positive and negative residuals, any quantile of interest can be modeled.

Similar to the linear regression model, quantile regression models assume that observations are independent of one another, and that some percentile of the response variable can be represented by a linear combination of a set of predictor variables. Perhaps the most advantageous aspect of quantile regression is the lack of assumptions regarding the distribution of the error.

While no assumptions about the error are required, estimation of standard errors for regression coefficients may be simplified by assuming the errors are identically and independently distributed (IID). When the IID assumption is not appropriate, bootstrap may be used to provide valid estimates of the standard errors for the regression coefficients.

CHAPTER 3

SIMULATION DESIGN

3.1 SCENARIO 1

The first simulation study compares SLR to QR, at the median only, when the normality of the errors assumption is violated, but variances are constant. For each scenario, the response variable was generated as $y_i = 15 + 0.4x_i + \varepsilon_i$ with the predictor, x_i , generated from a uniform distribution ranging from 0 to 2. A number of different sample sizes were included ($n=20, 30, 50, 100, 200, 300, 500$, and 650) in order to observe the performance of SLR and QR as sample size increased. To ensure a range of non-normal error distributions, data were generated with skewness levels of 0, 0.5, 1, 2, 2.5, 3, 4, and 6, while maintaining a standard deviation that was the same across all scenarios. Skewness was calculated using $\frac{\sum_j(z_j - \bar{z})^3/n}{s^3}$, where \bar{z} , s , and n represent the mean, standard deviation, and population size, respectively, and z_j denotes the j^{th} observation from the distribution. For all skewness levels, both left (negative) and right (positive) skewed data were generated.

With the exception of a skewness of 0, which utilized a normal distribution, the distribution from which these errors were selected were derived using gamma distributions with parameters that would yield the desired skewness while maintaining a consistent standard deviation. The probability distribution function of a Gamma distribution is defined as $f(u) = \left[\frac{1}{\Gamma(\alpha)\beta^\alpha} \right] u^{\alpha-1} e^{-u/\beta}$, where $\Gamma(\alpha) = \int_0^\infty u^{\alpha-1} e^{-u} du$, and

α and β represent the shape and scale parameters, respectively. The rate of a gamma distribution is the inverse of the scale parameter (Wackerly et. al. 2008). Once the appropriate error terms were found, the distributions were centered at their respective mean values. This ensured that, for all scenarios simulated, errors were pulled from distributions with means approximately zero and similar standard deviations of approximately 2.5. For simulations with negative skewness, the error distributions were identical, but y_i was generated by subtracting the error term (i.e $y_i = 15 + 0.4x_i - \varepsilon_i$).

The final error distributions are presented in Table 3.1. Histograms of these error distributions prior to centering are presented in Figures 3.1-3.7, along with a normal curve with identical means and standard deviations for comparison.

Simulations were designed to compare these two analytic methods based on measures related to β_1 , since inference for both SLR and QR are typically focused on β_1 , rather than the intercept β_0 . For each simulation scenario, 1,000 data sets were generated and the following measures calculated;

- 1) The average point estimate, defined as $\frac{1}{1000} \sum_{i=1}^{1000} \hat{\beta}_{1i}$
- 2) The average percent bias, defined as $\left(\frac{1}{1000} \sum_{i=1}^{1000} \frac{\hat{\beta}_{1i} - \beta_1}{\beta_1}\right) * 100$, which represents the average percent difference between the estimated $\hat{\beta}_1$ and the true value β_1 .
- 3) The 95% Confidence Interval was calculated for each estimate, defined as $\hat{\beta}_{1i} \pm t_{(0.975, n-2)} * \sqrt{var(\hat{\beta}_{1i})}$ where $t_{(0.975, n-2)}$ is the t value for the desired confidence level.

- 4) The average interval length was calculated as the average of the difference between the upper confidence limit (UCL) and the lower confidence limit (LCL), $\frac{1}{1000} \sum_{i=1}^{1000} (UCL_i - LCL_i)$.
- 5) Coverage Percentage was calculated as the percent of simulations for which the true value of the slope coefficient β_1 falls within the 95% confidence interval, $\frac{1}{1000} \sum_{i=1}^{1000} I(LCL_i < \beta_1 < UCL_i)$. We also created a confidence interval for the expected coverage percentage of 95%, defined as $0.95 \pm z_{0.975} * \frac{\sqrt{\frac{0.95}{1-0.95}}}{1000}$. This calculation resulted in an interval of (93.649, 96.351).
- 6) Power was calculated as the percent of simulations which were able to detect an association between the predictor and response variable. This was calculated as the percentage of simulations which had a significant slope estimate ($p - value < 0.05$), $\frac{1}{1000} \sum_{i=1}^{1000} I(p - value < 0.05)$.
- 7) The average percent magnitude of the difference between the predicted slope coefficient and the true slope, the mean absolute bias percent, was calculated as $\left(\frac{1}{1000} \sum_{i=1}^{1000} \frac{|\widehat{\beta}_{1i} - \beta_1|}{\beta_1} \right) * 100$.

Tables 3.2-3.9 present the results of the above calculations in cases where the error is positively skewed, and Tables 3.10-3.16 present results for negatively skewed errors. For all simulations, standard errors for quantile regression were generated via bootstrap with 200 replicates.

In order to verify that the tests have the proper type I error level, the same simulations were run, with $n=20, 200, \text{ and } 650$ and $y_i = 15 + 0 * x_i + \varepsilon_i$ (no association

between the predictor and the response). The significance level was obtained the same way that power was defined previously, and is presented in Table 3.17 by type of regression method used.

3.2 SCENARIO 2

The second simulation study compares the performance of linear regression to quantile regression at various percentiles of the response variable (10th, 25th, 50th, 75th, 90th), while violating the assumption of homoscedastic errors. Similar to the previous scenario, 1,000 data sets were generated at each sample size from $n=20, 50, 100, 200, 300, 500,$ and 650 . The generated response variable was defined as $y_i = 15 + 0.4x_i + \varepsilon_i * g(x_i)$ with x_i generated from a Uniform distribution ranging from 1 to 6, ε_i generated from a Normal distribution with mean 0 and standard deviation 2.3, and $g(x_i)$ representing a function of x_i . Multiplying the random error term by $g(x_i)$ yields a response variable that is heteroscedastic for the range of x . The specific function, $g(x_i)$, was chosen so that the error was a linear function of x and the slope relating the predictor to the 90th percentile of the outcome was approximately 1.5, 2, and 3 times that of the slope at the median. The final equations showing the error distributions used for this scenario are shown in Table 3.18.

This simulation scenario utilized the same methods of comparison between linear regression and quantile regression as was used in the previous simulation scenario, specifically coverage percentage and statistical power. Calculating the coverage percentage for different percentiles of quantile regression, however, required small alterations to the calculations. As the distribution of the error is known, the confidence intervals and coverage percentage for the quantile regression analyses were calculated

using the theoretical distribution to determine the true slope at a given percentile. As the violation of the homoscedasticity of errors suggests that the impact of the predictor upon the response variable may vary across percentiles, examining the difference between the known association at the mean and the estimated association at the quantiles is of interest.

The standard errors for the quantile regression analyses were again calculated using the bootstrap method with 200 iterations. The SLR results were compared to QR results at the 10th, 25th, 50th, 75th, and 90th percentiles. Tables 3.19-3.21 present the results of the comparisons between parameter estimates, and Tables 3.22-3.24 show the coverage percentages for each of our experiments.

The same simulations with $y_i = 15 + 0 * x_i + \varepsilon_i * g(x_i)$ (no association between the predictor and the response) for the same sample sizes, $n=20, 50, 100, 200, 300, 500, 650$ were conducted. This allowed for the verification that the SLR analyses maintained the proper type I error level. The magnitude of the estimates, as well as the statistical power in SLR versus quantile regression at the percentiles of interest were calculated, and these results are presented in Tables 3.25-3.27

3.3 SCENARIO 3

The final simulation scenario was designed to represent a biological example where the assumptions of both normality and homoscedasticity of errors would be unrealistic. This scenario allows a comparison of difference in conclusions and interpretation of results between these two analytic techniques. For this purpose, we chose to examine C-reactive protein (CRP), a non-negative marker of inflammation with a highly skewed distribution. Previous research using the Third National Health and

Nutrition Examination Survey (NHANES III), has shown mean CRP levels of approximately 4.3 mg/L and median levels of about 2.1 mg/L in men of all races without coronary heart disease aged 30-74 (Wong et. al., 2001). Similarly, a study of men without heart disease, hypertension, diabetes, cancer, asthma, or bronchitis aged 49-97 showed a similar baseline mean and median CRP levels of 3.8 mg/L and 2.3 mg/L, respectively (Bind, et. al. 2016).

As a marker of inflammation within the body, CRP has been found to increase as an acute response to vigorous physical activity. Weight, et. al (1991) compared baseline CRP levels from male and female competitive distance runners prior to beginning a marathon to levels immediately, 24 hours, 48 hours, and 6 days post completion. CRP levels in this sample peaked at 24 hours, decreased from this peak at 48 hours, and returned to baseline levels by the 6th day. Among male triathlon competitors, a similar pattern of results was found, with a peak in CRP levels at 24 hours after cessation of exercise, and returning to baseline after 48 hours (Taylor, et. al. 1987). Despite this acute response, several studies have been performed to understand the relationship between regular physical activity and baseline CRP level, and have largely found that those with a sedentary lifestyle have higher baseline CRP levels as compared with those who are more regularly physically active. Abramson et. al. (2002) found that, of adult men and women over 40, as frequency of physical activity within the past month increased, the odds of having elevated CRP levels decreased. Similarly, a case-control study found that baseline CRP levels in long-distance runners was significantly lower than CRP levels among an untrained control group (Tomaszewski, et. al 2003).

For the purposes of this experiment, it is assumed that the relationship between the predictor (physical activity level) and the outcome (CRP level) differs across the conditional distribution of CRP, and therefore by percentile, at least in part due to the time of measurement. For example, those within the 10th percentile of CRP level could have been measured more than 48 hours after physical activity, while those within the 90th percentile could have been measured within the peak period of about 24 hours post exercise. This assumption yields a scenario in which both the homoscedasticity and normality of errors assumptions are violated.

Bind et. al. (2016) conducted a study which enrolled 1,112 men aged 49-100 with no heart disease, hypertension, diabetes, cancer, recurrent asthma, or bronchitis. In this study, the mean, 5th percentile, median, and 95th percentile of baseline CRP levels were measured, yielding results of 3.8, 0.4, 2.3, and 24.5 mg/L, respectively. For this simulation scenario, the error was generated so that the error term was a linear function of the predictor such that the outcome variable of CRP would follow this general pattern with heteroscedasticity over the range of x . In order to mimic this scenario, x_i was generated from a random uniform distribution ranging from -2 to 2, with $y_i = 1 + 0.5 * x_i + \varepsilon_i * g(x_i)$. In this scenario, the error term was generated by multiplying $g(x_i) = \left(1 + 5 * x_i / 11\right)$ by a random gamma distribution with shape parameter of 0.15 and scale parameter of 18. As the goal is to create different variance in the outcome that was a function of the predictor, histograms of this error distribution are included at a range of values of the predictor in Figure 3.8-3.10. The y_i variable was then centered to ensure a mean value of 0. In this experiment, 1000 data sets were generated with a sample size of 50.

To compare the results from linear regression to quantile regression, identical methods as were discussed previously, with appropriate modifications to the calculations at the 10th, 25th, 50th, 75th, and 90th percentiles. The results from this simulation experiment are presented in Table 3.28.

Table 3.1: Simulation Scenario 1 Error Distributions by Skewness

Skewness	Calculated Skewness	Error Distribution	Std. Deviation
0	0.0007	Normal(mean=0, sd=2.5)	2.5092
0.5	0.5052	Gamma(shape=16, rate=1.6)	2.4916
1	1.0184	Gamma(shape=4, rate=0.8)	2.4947
2	2.0322	Gamma(shape=1, rate=0.4)	2.5066
2.5	2.5212	Gamma(shape=0.65, rate=0.32)	2.5260
3	3.0289	Gamma(shape=0.45, rate=0.27)	2.4920
4	4.0417	Gamma(shape=0.26, rate=0.2)	2.5694
6	5.9522	Gamma(shape=0.12, rate=0.14)	2.5310

Table 3.2: Estimates for Mean(SLR) and Median(QR) when Linear Regression Assumptions Hold ($\epsilon \sim N(0, sd = 2.5)$)

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4059	0.4571	0.4173	0.4381	0.4110	0.4016	0.3957	0.3814	0.3934	0.4052	0.3940	0.3863	0.4009	0.3903
Bias Percentage	1.4859	14.2713	4.3238	9.5273	2.7426	0.3993	-1.0739	-4.6417	-1.6530	1.3027	-1.5099	-3.4346	0.2211	-2.4317
Coverage Percent	96.0	96.9	94.3	96.0	94.8	96.2	94.3	94.0	94.9	94.8	95.5	93.3	94.4	95.2
Avg. Interval Length	4.2340	6.1550	2.4922	3.4547	1.7274	2.2896	1.2155	1.5995	0.9875	1.2808	0.7623	0.9753	0.6684	0.8571
Power	0.058	0.040	0.101	0.061	0.161	0.097	0.268	0.152	0.333	0.242	0.526	0.354	0.654	0.440
Mean Absolute Bias (%)	197.15	258.41	127.17	157.41	90.91	107.77	64.35	76.61	48.25	61.90	39.45	50.26	33.97	43.09

16

Table 3.3: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 0.5

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3538	0.3871	0.3936	0.3878	0.3914	0.4234	0.4032	0.4096	0.3921	0.3939	0.4029	0.3977	0.4018	0.4052
Bias Percentage	-11.556	-3.2341	-1.5952	-3.0484	-2.1387	5.8375	0.8077	2.4009	-1.9644	-1.5127	0.7325	-0.5821	0.4426	1.2878
Coverage Percent	94.0	97.1	96.2	96.5	94.4	94.5	94.2	95.4	95.2	93.8	95.3	94.5	95.6	93.5
Avg. Interval Length	4.1997	6.1471	2.4683	3.3368	1.7235	2.2620	1.2085	1.5555	0.9860	1.2680	0.7619	0.9658	0.6684	0.8458
Power	0.079	0.033	0.083	0.067	0.146	0.104	0.282	0.175	0.349	0.249	0.560	0.357	0.660	0.476
Mean Absolute Bias (%)	205.86	257.47	118.86	151.96	87.36	110.36	63.24	73.13	50.73	63.58	38.74	46.81	33.62	43.02

Table 3.4: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 1

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4365	0.3977	0.4269	0.4131	0.4173	0.4267	0.4135	0.4037	0.4012	0.4163	0.4075	0.4046	0.4041	0.4084
Bias Percentage	9.1359	-0.5632	6.7339	3.2782	4.3200	6.6829	3.3739	0.9280	0.3060	4.0744	1.8780	1.1573	1.0221	2.1108
Coverage Percent	94.7	97.7	95.3	95.9	96.1	96.0	96.3	94.8	95.3	95.8	95.6	94.0	94.4	94.8
Avg. Interval Length	4.1416	5.9296	2.4736	3.1961	1.7178	2.1929	1.2101	1.5024	0.9833	1.2312	0.7616	0.9309	0.6682	0.8191
Power	0.089	0.025	0.104	0.086	0.146	0.118	0.270	0.188	0.350	0.273	0.577	0.402	0.657	0.504
Mean Absolute Bias (%)	197.06	241.97	130.70	146.67	83.06	101.41	59.60	72.11	48.68	60.95	39.36	45.48	34.25	42.18

17

Table 3.5: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 2

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4003	0.3649	0.4134	0.4065	0.3849	0.3736	0.4059	0.3905	0.4104	0.4026	0.3907	0.3859	0.4087	0.4076
Bias Percentage	0.0848	-8.7812	3.3599	1.6354	-3.7650	-6.6023	1.4850	-2.3779	2.6058	0.6590	-2.3214	-3.5281	2.1760	1.8934
Coverage Percent	96.3	98.2	95.9	96.2	95.0	95.6	94.8	94.6	93.9	94.6	95.0	95.1	95.9	93.2
Avg. Interval Length	4.0660	5.0126	2.4759	2.6852	1.7236	1.8302	1.2060	1.2687	0.9813	1.0136	0.7616	0.7821	0.6671	0.6762
Power	0.069	0.033	0.101	0.098	0.131	0.116	0.291	0.237	0.362	0.366	0.527	0.502	0.681	0.644
Mean Absolute Bias (%)	188.47	197.87	120.67	123.06	86.06	83.51	62.12	61.61	50.45	50.72	39.82	39.11	31.83	34.05

Table 3.6: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 2.5

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4292	0.3954	0.3983	0.3962	0.4143	0.4063	0.3750	0.4111	0.4071	0.4072	0.3948	0.4045	0.3903	0.3965
Bias Percentage	7.3116	-1.1625	-0.4127	-0.9516	3.5691	1.5704	-6.2379	2.7766	1.7796	1.7918	-1.2896	1.1286	-2.4200	-0.8854
Coverage Percent	96.6	98.7	95.3	95.6	94.7	95.4	94.4	95.1	96.3	95.6	95.4	93.9	95.8	96.1
Avg. Interval Length	4.0328	4.5423	2.4708	2.3583	1.7239	1.5901	1.2140	1.1022	0.9921	0.8907	0.7699	0.6807	0.6720	0.5995
Power	0.071	0.026	0.109	0.092	0.168	0.177	0.241	0.318	0.389	0.440	0.524	0.635	0.617	0.733
Mean Absolute Bias (%)	191.98	168.81	125.05	106.29	88.03	73.18	61.79	54.06	51.25	41.79	37.16	34.84	34.37	28.15

18

Table 3.7: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 3

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3544	0.4040	0.3979	0.3933	0.4096	0.4022	0.4005	0.4040	0.3943	0.4049	0.4064	0.4046	0.3994	0.4034
Bias Percentage	-11.4108	1.0042	-0.5259	-1.6827	2.4072	0.5405	0.1174	1.0123	-1.4323	1.2187	1.6068	1.1447	-0.1491	0.8558
Coverage Percent	95.2	98.9	94.8	96.6	96.3	95.1	94.8	94.8	94.7	95.1	95.3	94.4	95.2	95.0
Avg. Interval Length	3.9478	4.0306	2.4235	1.9351	1.6950	1.2919	1.1994	0.8787	0.9755	0.7066	0.7551	0.5451	0.6637	0.4731
Power	0.073	0.038	0.121	0.119	0.154	0.252	0.266	0.455	0.362	0.620	0.556	0.816	0.679	0.913
Mean Absolute Bias (%)	187.21	142.42	120.53	86.44	81.77	61.50	61.37	42.84	50.83	35.07	37.89	27.01	32.60	23.15

Table 3.8: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 4

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4026	0.4065	0.3713	0.3960	0.3617	0.3945	0.3901	0.3980	0.4005	0.3952	0.4033	0.3962	0.3941	0.3972
Bias Percentage	0.6377	1.6234	-7.1776	-1.0011	-9.5814	-1.3684	-2.4630	-0.5077	0.1149	-1.2113	0.8364	-0.9422	-1.4867	-0.7038
Coverage Percent	96.0	99.4	94.7	98.8	95.5	97.5	95.4	95.3	94.6	95.2	94.4	95.6	95.1	94.8
Avg. Interval Length	3.9209	3.0002	2.4542	1.2411	1.7166	0.7532	1.2331	0.5023	0.9927	0.4054	0.7717	0.3104	0.6799	0.2705
Power	0.105	0.074	0.122	0.309	0.145	0.573	0.261	0.848	0.379	0.952	0.528	0.995	0.619	1.000
Mean Absolute Bias (%)	188.17	83.86	122.40	49.28	86.96	33.81	63.61	24.00	51.43	19.31	39.53	15.06	35.41	13.59

19

Table 3.9: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to 6

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3869	0.4003	0.3877	0.4006	0.4006	0.3982	0.4240	0.4003	0.3981	0.3998	0.4111	0.4000	0.4053	0.3997
Bias Percentage	-3.2726	0.0736	-3.0810	0.1492	0.1605	-0.4441	5.9917	0.0655	-0.4806	-0.0579	2.7680	-0.0023	1.3167	-0.0697
Coverage Percent	96.6	100.0	96.5	99.8	95.5	99.5	96.6	99.3	95.3	98.9	93.6	98.0	95.5	94.9
Avg. Interval Length	3.5649	1.7384	2.2696	0.3463	1.6392	0.1426	1.1736	0.0747	0.9620	0.0544	0.7460	0.0389	0.6605	0.0336
Power	0.155	0.338	0.155	0.930	0.208	0.998	0.347	1.000	0.404	1.000	0.592	1.000	0.682	1.000
Mean Absolute Bias (%)	174.46	18.53	115.56	7.05	85.12	4.38	59.70	2.79	48.87	2.30	39.30	1.77	32.01	1.61

Table 3.10: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -0.5

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4462	0.4129	0.4064	0.4122	0.4086	0.3767	0.3968	0.3904	0.4079	0.4061	0.3971	0.4023	0.3982	0.3948
Bias Percentage	11.5564	3.2341	1.5952	3.0484	2.1387	-5.8375	-0.8077	-2.4009	1.9644	1.5127	-0.7325	0.5821	-0.4426	-1.2878
Coverage Percent	94.0	97.1	96.2	96.5	94.4	94.6	94.2	95.9	95.2	93.7	95.3	94.5	95.6	93.7
Avg. Interval Length	4.1997	6.1470	2.4683	3.3368	1.7235	2.2620	1.2085	1.5555	0.9860	1.2680	0.7619	0.9658	0.6684	0.8457
Power	0.076	0.033	0.097	0.073	0.143	0.119	0.267	0.175	0.379	0.243	0.519	0.377	0.655	0.474
Mean Absolute Bias (%)	205.86	257.47	118.86	151.96	87.36	110.36	63.24	73.13	50.73	63.58	38.74	46.81	33.62	43.02

20

Table 3.11: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -1

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3635	0.4023	0.3731	0.3869	0.3827	0.3733	0.3865	0.3963	0.3988	0.3837	0.3925	0.3954	0.3959	0.3916
Bias Percentage	-9.1359	0.5632	-6.7339	-3.2782	-4.320	-6.6829	-3.3739	-0.9280	-0.3060	-4.0744	-1.8780	-1.1573	-1.0221	-2.1108
Coverage Percent	94.7	97.5	95.3	95.8	96.1	96.0	96.3	94.8	95.3	95.8	95.6	94.6	94.4	94.4
Avg. Interval Length	4.1416	5.9296	2.4736	3.1961	1.7178	2.1929	1.2101	1.5024	0.9833	1.2312	0.7616	0.9309	0.6682	0.8191
Power	0.067	0.039	0.108	0.073	0.133	0.099	0.253	0.180	0.348	0.237	0.515	0.386	0.639	0.479
Mean Absolute Bias (%)	194.06	241.97	130.70	146.67	83.06	101.41	59.60	72.11	48.68	60.95	39.36	45.48	34.25	42.18

Table 3.12: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -2

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3997	0.4351	0.3866	0.3935	0.4151	0.4264	0.3941	0.4095	0.3896	0.3974	0.4093	0.4141	0.3913	0.3924
Bias Percentage	-0.0848	8.7812	-3.3599	-1.6354	3.7650	6.6023	-1.4850	2.3779	-2.6058	-0.6590	2.3214	3.5281	-2.1760	-1.8934
Coverage Percent	96.3	97.9	95.9	96.4	95.0	95.2	94.8	95.2	93.9	94.5	95.0	94.9	95.9	93.3
Avg. Interval Length	4.0660	5.0126	2.4759	2.6852	1.7236	1.8301	1.2060	1.2687	0.9813	1.0136	0.7616	0.7821	0.6671	0.6762
Power	0.058	0.036	0.080	0.093	0.166	0.149	0.248	0.251	0.345	0.350	0.566	0.562	0.645	0.620
Mean Absolute Bias (%)	188.47	197.87	120.67	123.06	86.06	83.51	62.12	61.61	50.45	50.72	39.82	39.11	31.83	34.05

21

Table 3.13: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -2.5

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3708	0.4046	0.4017	0.4038	0.3857	0.3937	0.4250	0.3889	0.3929	0.3928	0.4052	0.3955	0.4097	0.4035
Bias Percentage	-7.3116	1.1625	0.4127	0.9516	-3.5691	-1.5704	6.2379	-2.7766	-1.7796	-1.7918	1.2896	-1.1286	2.4200	0.8854
Coverage Percent	96.6	98.9	95.3	95.5	94.7	95.6	94.4	95.5	96.3	95.3	95.4	94.3	95.8	95.7
Avg. Interval Length	4.0328	4.5423	2.4708	2.3582	1.7239	1.5901	1.2140	1.1022	0.9921	0.8907	0.7699	0.6807	0.6720	0.5995
Power	0.062	0.025	0.120	0.120	0.148	0.154	0.281	0.289	0.357	0.423	0.556	0.607	0.663	0.758
Mean Absolute Bias (%)	191.98	168.81	125.05	106.29	88.03	73.18	61.79	54.06	51.25	41.79	37.15	34.84	34.37	28.15

Table 3.14: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -3

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4456	0.3960	0.4021	0.4067	0.3904	0.3978	0.3995	0.3960	0.4057	0.3951	0.3936	0.3954	0.4006	0.3966
Bias Percentage	11.411	-1.0042	0.5259	1.6827	-2.4072	-0.5405	-0.1174	-1.0123	1.4323	-1.2187	-1.6068	-1.1447	0.1491	-0.8558
Coverage Percent	95.2	98.9	94.8	96.8	96.3	95.3	94.8	94.7	94.7	95.0	95.3	94.4	95.2	94.9
Avg. Interval Length	3.9478	4.0306	2.4235	1.9351	1.6950	1.2919	1.1994	0.8786	0.9755	0.7066	0.7551	0.5451	0.6637	0.4731
Power	0.076	0.045	0.117	0.145	0.167	0.257	0.271	0.430	0.383	0.577	0.558	0.810	0.668	0.902
Mean Absolute Bias (%)	187.21	142.42	120.53	86.44	81.77	61.50	61.37	42.84	50.83	35.07	37.89	27.01	32.60	23.15

22

Table 3.15: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -4

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.3975	0.3935	0.4287	0.4040	0.4383	0.4055	0.4099	0.4020	0.3995	0.4048	0.3967	0.4038	0.4059	0.4028
Bias Percentage	-0.6377	-1.6234	7.1776	1.0011	9.5814	1.3684	2.4630	0.5077	-0.1149	1.2113	-0.8364	0.9422	1.4867	0.7038
Coverage Percent	96.0	99.3	94.7	98.8	95.5	97.5	95.4	96.3	94.6	95.1	94.4	94.9	95.1	95.0
Avg. Interval Length	3.9209	3.0002	2.4542	1.2411	1.7166	0.7532	1.2331	0.5023	0.9927	0.4054	0.7717	0.3104	0.6799	0.2705
Power	0.080	0.062	0.129	0.316	0.187	0.619	0.271	0.869	0.378	0.956	0.553	0.995	0.646	0.998
Mean Absolute Bias (%)	188.17	83.86	122.40	49.28	86.96	33.81	63.61	24.00	51.43	19.31	39.53	15.06	35.41	13.59

Table 3.16: Estimates for Mean(SLR) and Median(QR) with Error Distribution Skewed to -6

Sample Size	20		50		100		200		300		500		650	
Method	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR	SLR	QR
Avg. estimate	0.4131	0.3997	0.4123	0.3994	0.3994	0.4018	0.3760	0.3997	0.4019	0.4002	0.3889	0.4000	0.3947	0.4003
Bias Percentage	3.2726	-0.0736	3.0810	-0.1492	-0.1605	0.4441	-5.9917	-0.0655	0.4806	0.0579	-2.7680	0.0023	-1.3167	0.0697
Coverage Percent	96.6	100	96.5	99.8	95.5	99.7	96.6	99.4	95.3	98.8	93.6	98.2	95.5	94.8
Avg. Interval Length	3.5649	1.7382	2.2696	0.3463	1.6392	0.1426	1.1736	0.0747	0.9620	0.0542	0.7460	0.0389	0.6605	0.0336
Power	0.143	0.330	0.171	0.926	0.216	0.998	0.273	1.000	0.380	1.000	0.535	1.000	0.658	1.000
Mean Absolute Bias (%)	174.46	18.53	115.56	7.05	86.12	4.38	59.70	2.79	48.87	2.30	39.30	1.77	32.01	1.61

Table 3.17: Type I Error Rate for SLR and Median QR with Slope Coefficient of 0

Skewness of Error	Sample Size	Type I Error Rate of SLR	Type I Error Rate of QR
0	20	0.040	0.035
	200	0.057	0.058
	650	0.056	0.048
0.5	20	0.060	0.029
	200	0.058	0.045
	650	0.044	0.064
1	20	0.053	0.022
	200	0.037	0.052
	650	0.056	0.054
2	20	0.037	0.019
	200	0.052	0.055
	650	0.041	0.070
2.5	20	0.034	0.009
	200	0.056	0.047
	650	0.042	0.043
3	20	0.048	0.012
	200	0.052	0.053
	650	0.048	0.051
4	20	0.040	0.007
	200	0.046	0.041
	650	0.049	0.054
6	20	0.034	0.000
	200	0.034	0.007
	650	0.045	0.047

Table 3.18: Simulation Scenario 2 Error Distributions

$\beta_1^{(90)}: \beta_1^{(50)}$	Error Distribution
1	$Normal(0,2.3)$
1.5	$\left(1 + \left(\frac{x}{14}\right)\right) * Normal(0,2.3)$
2	$\left(1 + \left(\frac{x}{7}\right)\right) * Normal(0,2.3)$
3	$\left(1 + \frac{3 * x}{11}\right) * Normal(0,2.3)$

Table 3.19: Comparison of SLR(Mean) and QR(p th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5

Sample Size	QR (10 th) True Slope: 0.1896		QR (25 th) True Slope: 0.2892		QR (50 th) True Slope: 0.4000		QR (75 th) True Slope: 0.5108		QR (90 th) True Slope: 0.6105		Mean
	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate
20	0.1615	-59.63	0.2787	-30.33	0.4251	6.28	0.4897	22.43	0.5911	47.78	0.4004
50	0.1828	-54.30	0.2844	-28.90	0.4180	4.50	0.5196	29.90	0.5976	49.40	0.4062
100	0.1696	-57.60	0.2875	-28.13	0.4009	0.23	0.5172	29.30	0.6235	55.88	0.4054
200	0.1838	-54.05	0.2820	-29.50	0.3908	-2.30	0.5006	25.15	0.6131	53.28	0.3975
300	0.1872	-53.20	0.2861	-28.48	0.4018	0.45	0.5147	28.68	0.6187	54.68	0.3968
500	0.1896	-52.60	0.2843	-28.93	0.3931	-1.73	0.5021	25.53	0.6034	50.85	0.3970
650	0.1859	-53.53	0.2869	-28.28	0.3956	-1.10	0.5048	26.20	0.6070	51.75	0.4002

25

Table 3.20: Comparison of SLR(Mean) and QR(p th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2

Sample Size	QR (10 th) True Slope: -0.0211		QR (25 th) True Slope: 0.1784		QR (50 th) True Slope: 0.4000		QR (75 th) True Slope: 0.6216		QR (90 th) True Slope: 0.8211		Mean
	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate
20	-0.0386	-109.65	0.1695	-57.63	0.4282	7.05	0.5898	47.45	0.7801	95.03	0.3986
50	-0.0226	-105.65	0.1748	-56.30	0.4216	5.40	0.6317	57.93	0.8011	100.28	0.4061
100	-0.0418	-110.45	0.1778	-55.55	0.4008	0.20	0.6292	57.30	0.8328	108.20	0.4067
200	-0.0260	-106.50	0.1693	-57.68	0.3881	-2.98	0.6081	52.03	0.8217	105.43	0.3965
300	-0.0234	-105.85	0.1744	-56.40	0.4018	0.45	0.6259	56.48	0.8294	107.35	0.3961
500	-0.0205	-105.13	0.1727	-56.83	0.3916	-2.10	0.6112	52.80	0.8119	102.98	0.3962
650	-0.0250	-106.25	0.1762	-55.95	0.3949	-1.28	0.6145	53.63	0.8169	104.23	0.4000

Table 3.21: Comparison of SLR(Mean) and QR(p th percentile) Estimates in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3

Sample Size	QR (10 th) True Slope: -0.4039		QR (25 th) True Slope: -0.0231		QR (50 th) True Slope: 0.4000		QR (75 th) True Slope: 0.8231		QR (90 th) True Slope: 1.2039		Mean
	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate	Change from Mean (%)	Avg. Estimate
20	-0.4019	-200.48	-0.0283	-107.08	0.4340	8.50	0.7722	93.05	1.1222	180.55	0.3953
50	-0.3952	-198.80	-0.0236	-105.90	0.4277	6.93	0.8348	108.70	1.1706	192.65	0.4057
100	-0.4266	-206.65	-0.0221	-105.53	0.4003	0.08	0.8321	108.03	1.2120	203.00	0.4091
200	-0.4075	-201.88	-0.0356	-108.90	0.3838	-4.05	0.8041	101.03	1.2002	200.05	0.3948
300	-0.4055	-201.38	-0.0287	-107.18	0.4014	0.35	0.8280	107.00	1.2124	203.10	0.3947
500	-0.4020	-200.50	-0.0303	-107.58	0.3891	-2.73	0.8097	102.43	1.1912	197.80	0.3948
650	-0.4084	-202.10	-0.0249	-106.23	0.3935	-1.63	0.8144	103.60	1.1988	199.70	0.3997

26

Table 3.22: Comparison of SLR(Mean) and QR(p th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5

Sample Size	QR (10 th)	QR (25 th)	QR (50 th)	QR (75 th)	QR (90 th)	Mean
20	94.9	96.5	96.1	97.1	94.0	95.7
50	95.2	95.7	95.9	95.2	94.2	94.4
100	96.2	95.3	95.7	95.6	94.6	94.4
200	95.5	96.2	94.7	94.4	94.1	94.3
300	93.9	95.4	95.3	95.5	95.4	95.1
500	95.0	93.6	93.4	94.2	94.2	95.8
650	93.8	95.8	94.9	94.2	94.7	94.6

Table 3.23: Comparison of SLR(Mean) and QR(p th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2

Sample Size	QR (10 th)	QR (25 th)	QR (50 th)	QR (75 th)	QR (90 th)	Mean
20	94.7	96.5	96.0	94.2	93.8	96.1
50	95.5	95.9	96.2	95.6	94.0	94.6
100	96.4	95.0	95.7	95.0	95.0	94.5
200	95.4	96.1	94.9	93.9	94.2	94.4
300	93.8	95.6	95.5	95.2	95.8	95.0
500	95.4	93.5	93.2	94.8	94.0	95.7
650	94.7	95.9	94.9	94.2	94.8	94.2

Table 3.24: Comparison of SLR(Mean) and QR(p th percentile) Coverage Percentage in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0.4 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3

Sample Size	QR (10 th)	QR (25 th)	QR (50 th)	QR (75 th)	QR (90 th)	Mean
20	94.3	96.2	95.3	96.4	93.9	95.8
50	95.5	96.0	96.2	95.7	93.5	94.7
100	96.5	95.2	95.7	94.9	94.5	94.3
200	95.1	96.2	95.0	93.5	94.4	94.0
300	93.2	95.7	95.2	95.2	95.7	94.9
500	95.5	94.3	93.3	94.2	94.2	95.6
650	94.8	96.1	95.0	94.1	95.2	93.9

Table 3.25: Comparison of SLR(Mean) and QR(*p*th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 1.5

Sample Size	QR (10 th) True Slope: -0.2105		QR (25 th) True Slope: -0.1108		QR (50 th) True Slope: 0.0000		QR (75 th) True Slope: 0.1108		QR (90 th) True Slope: 0.2105		Mean	
	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power ¹
20	-0.2385	0.053	-0.1213	0.038	0.0251	0.041	0.0897	0.035	0.1911	0.068	0.0004	0.043
50	-0.2172	0.069	-0.1156	0.042	0.0180	0.035	0.1196	0.065	0.1976	0.078	0.0062	0.056
100	-0.2034	0.087	-0.1125	0.073	0.0009	0.040	0.1172	0.075	0.2235	0.088	0.0054	0.056
200	-0.2162	0.131	-0.1180	0.083	-0.0092	0.053	0.1006	0.083	0.2131	0.144	-0.0025	0.057
300	-0.2128	0.201	-0.1139	0.120	0.0018	0.042	0.1147	0.115	0.2187	0.180	-0.0032	0.049
500	-0.2104	0.283	-0.1157	0.155	-0.0069	0.063	0.1021	0.133	0.2034	0.258	-0.0030	0.042
650	-0.2141	0.341	-0.1131	0.162	-0.0044	0.050	0.1048	0.173	0.2070	0.319	0.0002	0.054

¹ Power for SLR when the mean effect is truly 0 is equivalent to the type I error rate.

28

Table 3.26: Comparison of SLR(Mean) and QR(*p*th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 2

Sample Size	QR (10 th) True Slope: -0.4211		QR (25 th) True Slope: -0.2216		QR (50 th) True Slope: 0.0000		QR (75 th) True Slope: 0.2216		QR (90 th) True Slope: 0.4211		Mean	
	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power ¹
20	-0.4386	0.070	-0.2305	0.040	0.0282	0.044	0.1898	0.041	0.3801	0.077	-0.0014	0.039
50	-0.4226	0.098	-0.2252	0.051	0.0216	0.036	0.2317	0.084	0.4011	0.119	0.0060	0.054
100	-0.4418	0.181	-0.2222	0.093	0.0008	0.042	0.2292	0.105	0.4328	0.152	0.0067	0.055
200	-0.4260	0.302	-0.2307	0.159	-0.0119	0.052	0.2081	0.155	0.4217	0.293	-0.0035	0.056
300	-0.4234	0.429	-0.2256	0.228	0.0018	0.044	0.2259	0.214	0.4294	0.409	-0.0039	0.050
500	-0.4205	0.608	-0.2273	0.340	-0.0084	0.063	0.2112	0.310	0.4119	0.589	-0.0038	0.043
650	-0.4250	0.727	-0.2238	0.384	-0.0051	0.050	0.2145	0.365	0.4169	0.697	-0.0000	0.058

¹ Power for SLR when the mean effect is truly 0 is equivalent to the type I error rate.

Table 3.27: Comparison of SLR(Mean) and QR(p th percentile) Estimates of Power in Cases Where Slope Estimate is Not Consistent Across Distribution of the Outcome Due to Heterogeneity in Error Variance, with Mean Effect 0 and Ratio of $\beta_1^{(90)}$ and $\beta_1^{(50)}$ of 3

Sample Size	QR (10 th) True Slope: -0.8039		QR (25 th) True Slope: -0.4231		QR (50 th) True Slope: 0.0000		QR (75 th) True Slope: 0.4231		QR (90 th) True Slope: 0.8039		Mean	
	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power	Avg. Estimate	Power ¹
20	-0.8019	0.101	-0.4283	0.052	0.0340	0.046	0.3722	0.051	0.7222	0.097	-0.0047	0.042
50	-0.7952	0.169	-0.4236	0.076	0.0277	0.038	0.4348	0.110	0.7706	0.176	0.0057	0.053
100	-0.8266	0.332	-0.4220	0.156	0.0003	0.043	0.4321	0.180	0.8120	0.304	0.0091	0.057
200	-0.8075	0.518	-0.4356	0.292	-0.0162	0.051	0.4041	0.258	0.8002	0.519	-0.0052	0.060
300	-0.8055	0.693	-0.4287	0.404	0.0014	0.045	0.4280	0.391	0.8124	0.724	-0.0053	0.051
500	-0.8020	0.899	-0.4303	0.597	-0.0109	0.059	0.4097	0.563	0.7912	0.889	-0.0052	0.044
650	-0.8084	0.958	-0.4249	0.719	-0.0065	0.050	0.4144	0.693	0.7988	0.959	-0.0003	0.061

¹ Power for SLR when the mean effect is truly 0 is equivalent to the type I error rate.

29

Table 3.28: Estimates for SLR and QR(p th percentile) with Skewed Errors and Heteroscedasticity

Sample Size	50					
Method	SLR	QR (10 th)	QR (25 th)	QR (50 th)	QR (75 th)	QR (90 th)
True Slope	0.5000	0.5000	0.5005	0.5510	1.3253	4.1393
Avg. estimate	1.6943	0.5001	0.5020	0.5742	1.3562	4.0309
Bias Percentage	238.8618	0.0137	0.2946	4.2063	2.3322	-2.6170
Coverage Percent	80.2	100	100	97.4	94.7	93.1
Avg. Interval Length	3.5085	0.0096	0.0757	0.7779	4.5100	14.4338
Power	0.522	1.000	0.999	0.848	0.178	0.188
Mean Absolute Bias (%)	244.70	0.02	0.58	15.27	51.16	56.13

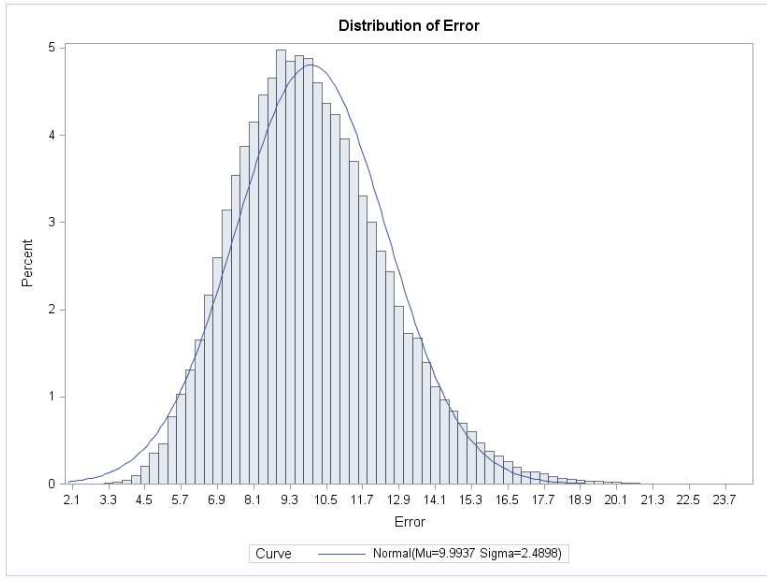


Figure 3.1: Scenario 1: Histogram of Error with Skewness 0.5 compared with Normal Curve

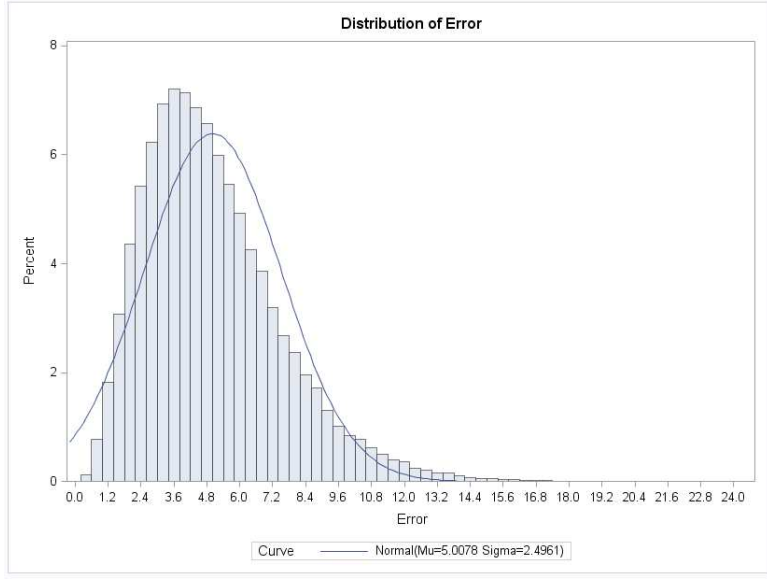


Figure 3.2: Scenario 1: Histogram of Error with Skewness 1 compared with Normal Curve

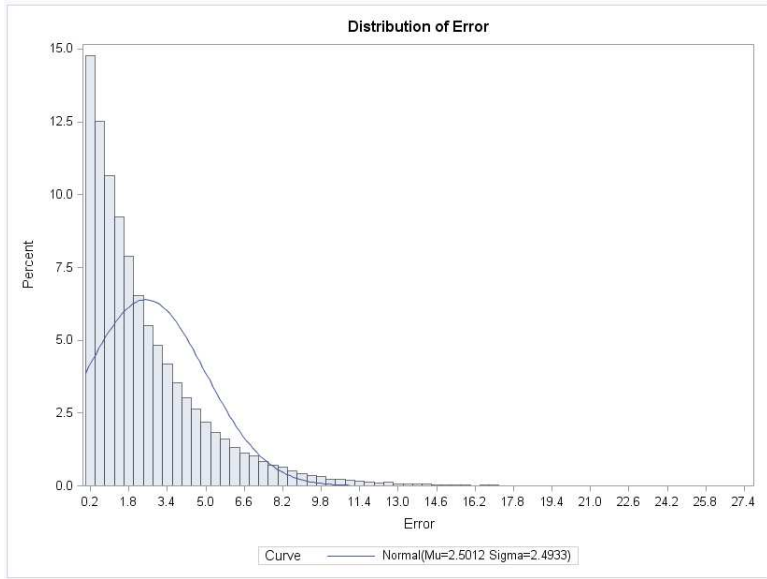


Figure 3.3: Scenario 1: Histogram of Error with Skewness 2 compared with Normal Curve

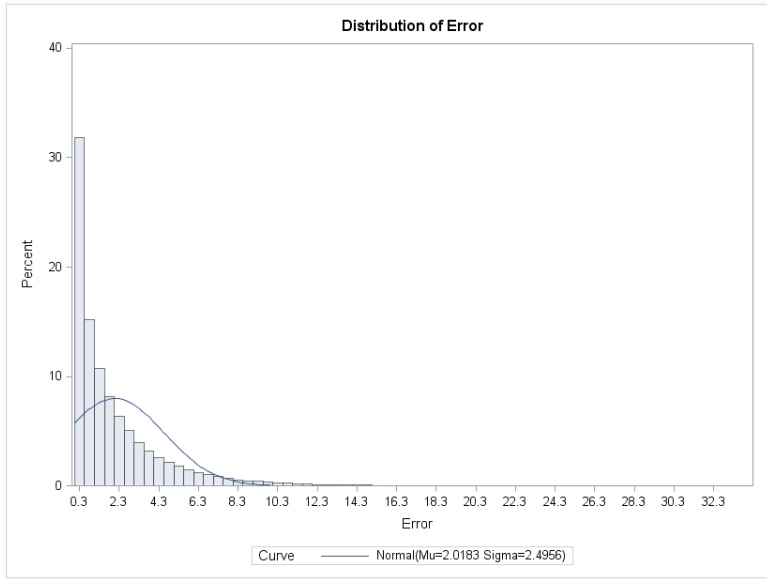


Figure 3.4: Scenario 1: Histogram of Error with Skewness 2.5 Compared with Normal Curve

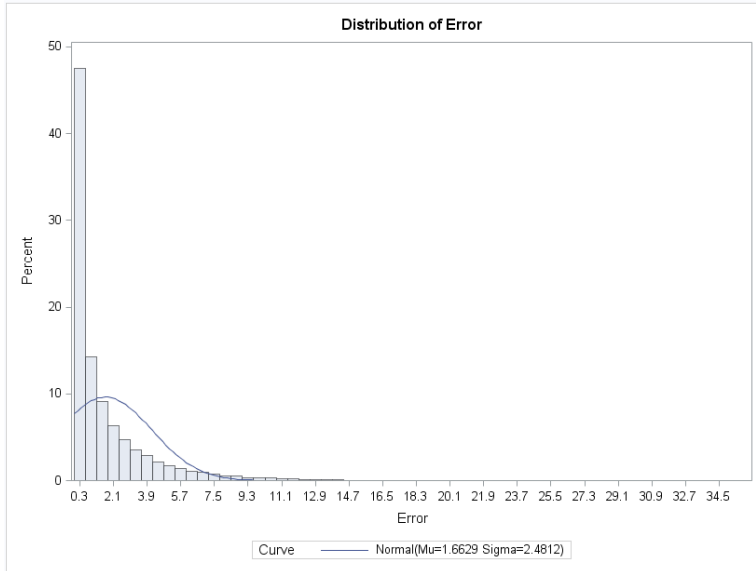


Figure 3.5: Scenario 1: Histogram of Error with Skewness 3 Compared with Normal Curve

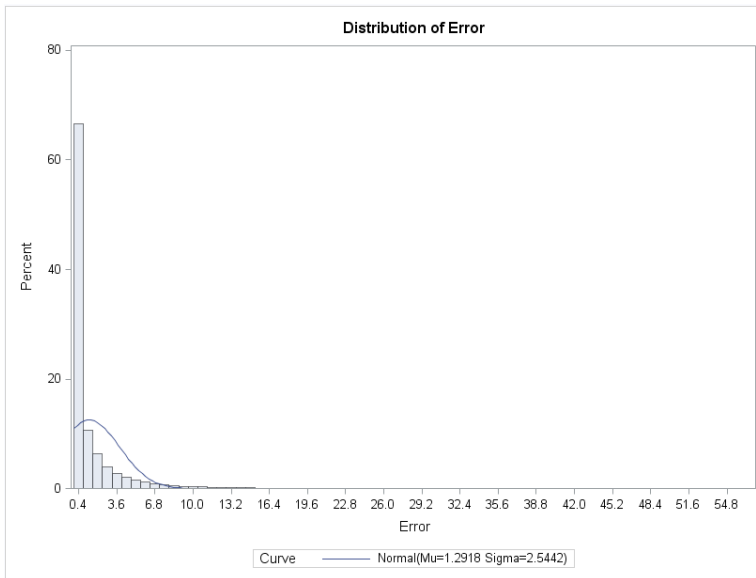


Figure 3.6: Scenario 1: Histogram of Error with Skewness 4 Compared with Normal Curve

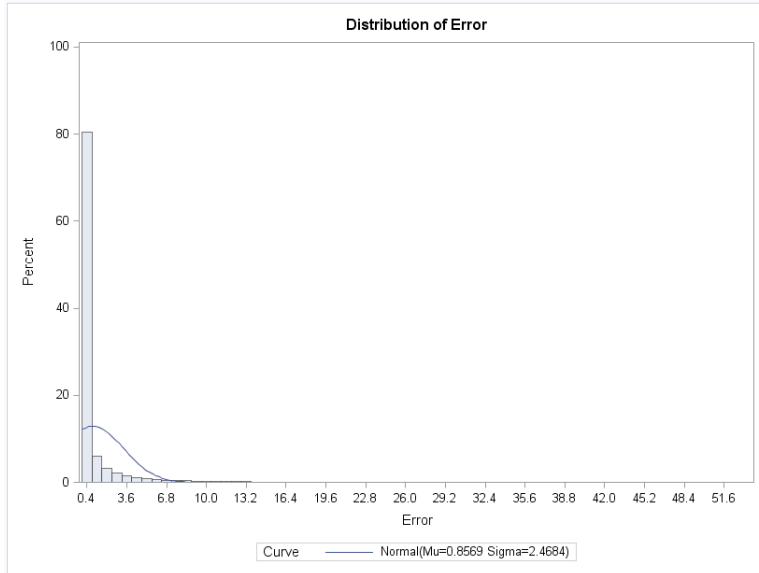


Figure 3.7: Scenario 1: Histogram of Error with Skewness 6 Compared with Normal Curve

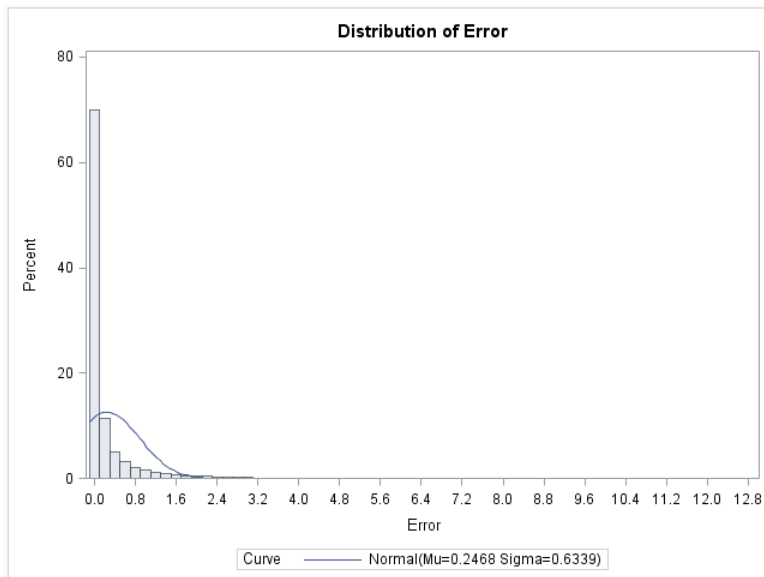


Figure 3.8: Scenario 3: Histogram of Error with Skewness and Heteroscedasticity at $x=-2$

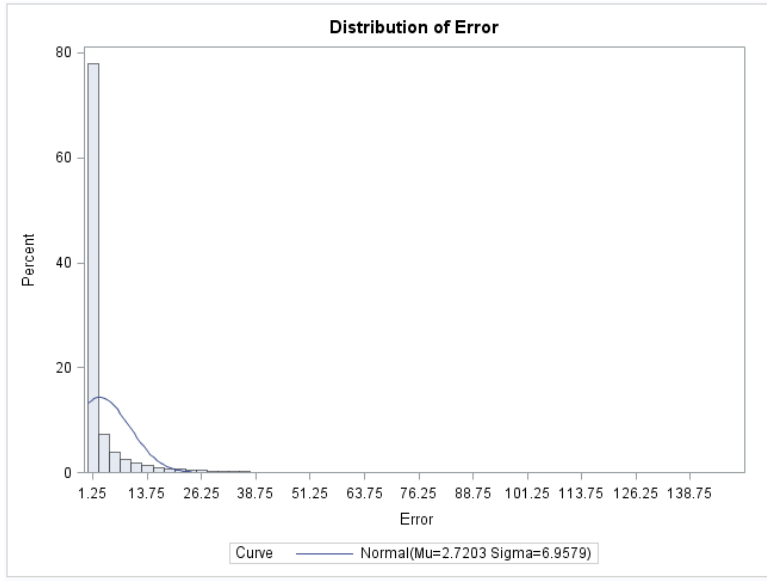


Figure 3.9: Scenario 3: Histogram of Error with Skewness and Heteroscedasticity at x=0

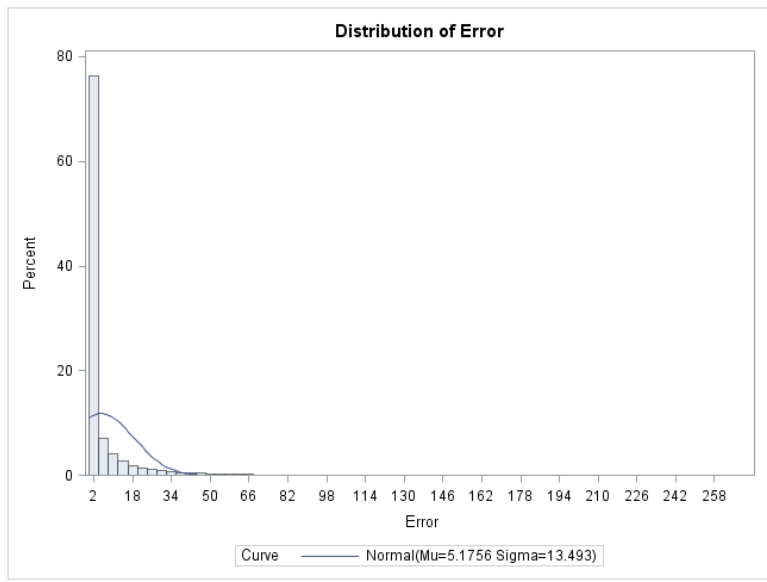


Figure 3.10: Scenario 3: Histogram of Error with Skewness and Heteroscedasticity at x=2

CHAPTER 4

RESULTS

4.1 SCENARIO 1

The results of the simulation are presented in Tables 3.2-3.16. The results in which the error distribution was negatively skewed presented very similar results to those from the positively skewed simulations. As such, the primary focus will be on the positively skewed results.

The coverage percentage represents the percentage of simulations for which the true slope was within the confidence interval for the estimated slope. Regardless of the level of skewness seen in the error distribution, the coverage percentage for the SLR models remain strong, suggesting that violations from the assumptions of normality of error do not result in any noticeable shifts in the coverage percentage. Figure 4.1 displays each of the SLR simulations' coverage percentage by skewness and sample size along with the calculated confidence interval of (93.649, 96.351) for the expected coverage percentage of 95%. The vast majority of the simulations fall within this 95% confidence limit, suggesting in most cases there is no evidence that violations of the normality assumption of random errors results in a change in coverage probability in SLR. Figure 4.2 represents the same graph of the coverage percentage and confidence bound for the QR models, also by skewness of the error distribution and sample size. For the majority of cases, the coverage percentage falls either within the calculated 95% confidence interval or above the upper confidence bound. In cases where the coverage is above the

upper confidence bound, this suggests that the confidence intervals for the parameter estimates, $\widehat{\beta}_{1l}$, are more conservative and thus more likely to include the true slope value β_1 .

In the case where the assumptions of linear regression hold, the SLR model was consistently more powerful than the QR model. Up to a skewness in the error term of about 2, linear regression is a more statistically powerful method of analysis than quantile regression. Once the skewness of the error surpasses this threshold, quantile regression has higher power to detect effects than linear regression. Prior to this threshold, the difference in power between SLR and QR is relatively small in smaller sample sizes, while there is a marked increase in difference of power level in small samples sizes after skewness exceeds the threshold. Figure 4.3 displays a comparison between statistical power of SLR versus QR for selected sample sizes. Similar patterns were observed among the excluded sample sizes, which were omitted to avoid clutter. It should be noted that, for a small sample size (e.g. $n=20, 50$), the statistical power of quantile regression remains lower than that of linear regression until the error distributions have a skewness greater than 3. This is likely due to the fact that, with a small sample size, each small sample is less likely to contain extreme error values selected from the larger population. Without these extreme values, the distribution of the error terms is impacted less by the non-normality of the error terms and linear regression remains a more powerful method of analysis.

Similar to the statistical power, the mean absolute bias percent of the models when the normality assumption is met is consistently smaller in the SLR model as compared with QR. As skewness of the error distribution increases, this difference

decreases until ultimately the QR models have a smaller mean absolute bias percentage. The difference initially becomes apparent among the larger sample sizes ($n=100$ and greater) when the skewness of the error is approximately 2. When the skewness of the error increases to 2.5 and greater, the difference becomes larger and more consistent for all sample sizes. Figure 4.4 shows a graphical representation of the difference in mean absolute bias percent between linear and quantile regression, by skewness and sample size.

Figure 4.5 displays the type I error rate for both SLR and QR at the median as a function of sample size and skewness of the error term when the slope coefficient is 0. With a Type I error rate of 5%, we would expect the simulations to report an association in 5% of the simulations, even though no association is present. This chart also includes a 95% confidence interval for this Type I error rate of 5% of (0.0365, 0.635). The linear regression simulations fall within this confidence interval as expected. The majority of the quantile regression simulations also fall within this interval, with the exception of the simulations run with a sample size of 20. For this sample size, the type I error rate was consistently below the confidence interval, suggesting that the type I error rate is actually less than the assumed 5%.

4.2 SCENARIO 2

This simulation experiment compares the linear regression estimate to the quantile regression estimates at various percentiles of interest, shown in Tables 3.19-3.21, as we increase the amount of variability seen in the distribution of the outcome as a factor of the predictor.

When the homoscedasticity of errors assumption is violated, the coverage percentage of the linear regression models remains strong regardless of the level of heteroscedasticity introduced to the model. Figure 4.6 shows the coverage percentage as a function of the sample size and heteroscedasticity factor for the linear regression models, while Figures 4.7 and 4.8 show the coverage percentage of the 10th and 90th percentiles, respectively. Coverage percentages in the quantile regression models not included as figures presented very similar results as the 10th and 90th percentile models.

The primary interest in this case is the difference between the estimates of the association between predictor and outcome. The homoscedasticity assumption of SLR implies that slope of the association between the predictor and the outcome remains the same at all percentiles of the conditional distribution of the outcome. That is, that the relationship between the predictor and the outcome at the mean of the conditional distribution is the same as that of any other percentile. Because the homoscedasticity assumption is violated in this scenario, the estimates of the slope at different percentiles of the outcome will be different from that of the slope at the mean. As soon as heteroscedasticity is introduced in our simulation experiments, the estimate from linear regression over or underestimates the association between the slope and response variables at our chosen percentiles of interest.

Considering the 10th percentile, even at our lowest level of dependence of variance upon the predictor, using the estimates from SLR to explain the association between the outcome and predictor in this percentile overestimates the true association by around 50%. As the strength of the dependence of variance upon the predictor gets

progressively stronger, this overestimation increases to 100% and 200%. At this percentile, the association between the predictor and response variables is truly an inverse relationship in our experiment, while the SLR estimate assumes that the relationship is positive. In the 90th percentile, similar patterns are seen in the magnitude of the differences in estimates, but the true relationship between predictor and response within the 90th percentile is underestimated by the SLR estimates. The underestimation climbs from approximately 50% to about 200%. Even within the 25th and 75th percentiles, the association is over and underpredicted, respectively, by the linear regression results.

Figure 4.9 displays the type I error rate for linear regression and shows that the type I error rate in our experiment remains within the 95% confidence interval for the expected rate of 5% regardless of the strength of the variance's dependence on the predictor. However, while the linear regression results indicate that there is no association between the predictor and outcome variables, these results are unable to accurately describe the true relationship for the conditional distribution of the outcome. That is, that there is a negative association within the lower percentiles and a positive association in the upper percentiles. Tables 3.25-3.27 display the estimates of the association and the statistical power of the quantile regression models at percentiles of interest when the slope at the mean is set to 0. Quantile regression is able to better capture the relationship at the upper and lower percentiles, as the statistical power is much higher than that of the corresponding linear regression models.

4.3 SCENARIO 3

The results from this simulation experiment are presented in Table 3.28. With the introduction of skewness and heteroscedasticity simultaneously existing within the error

term, inference and interpretation of the linear regression model are adversely affected. The linear regression estimates of the β_1 coefficient drastically overestimate the association for the majority of the percentiles of the conditional distributions, and drastically underestimate this relationship within the 90th percentile. On the other hand, the results from the quantile regression were able to more accurately estimate the relationship between the predictor variable and CRP. Considering all of the percentiles simultaneously gives a more complete picture of the relationship between the predictor and response variables. Coverage percentage of the SLR model dropped outside of what is expected with a type I error rate of 5%, suggesting that inferences made from linear regression in cases where both of these assumptions are violated might not be valid. The coverage remained steady in the quantile regression model, suggesting inference from quantile regression is valid across quantiles.

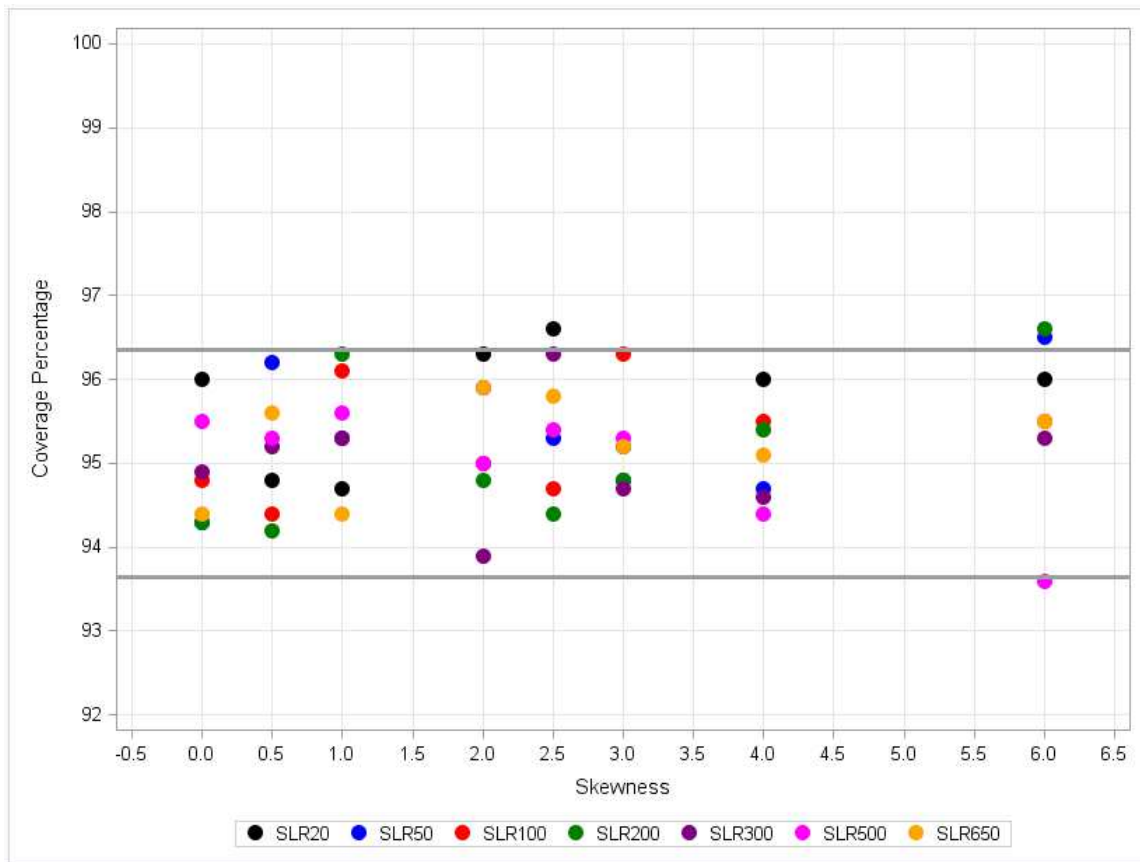


Figure 4.1: Coverage Percentage by Skewness and Sample size for SLR (Slope =0.4)

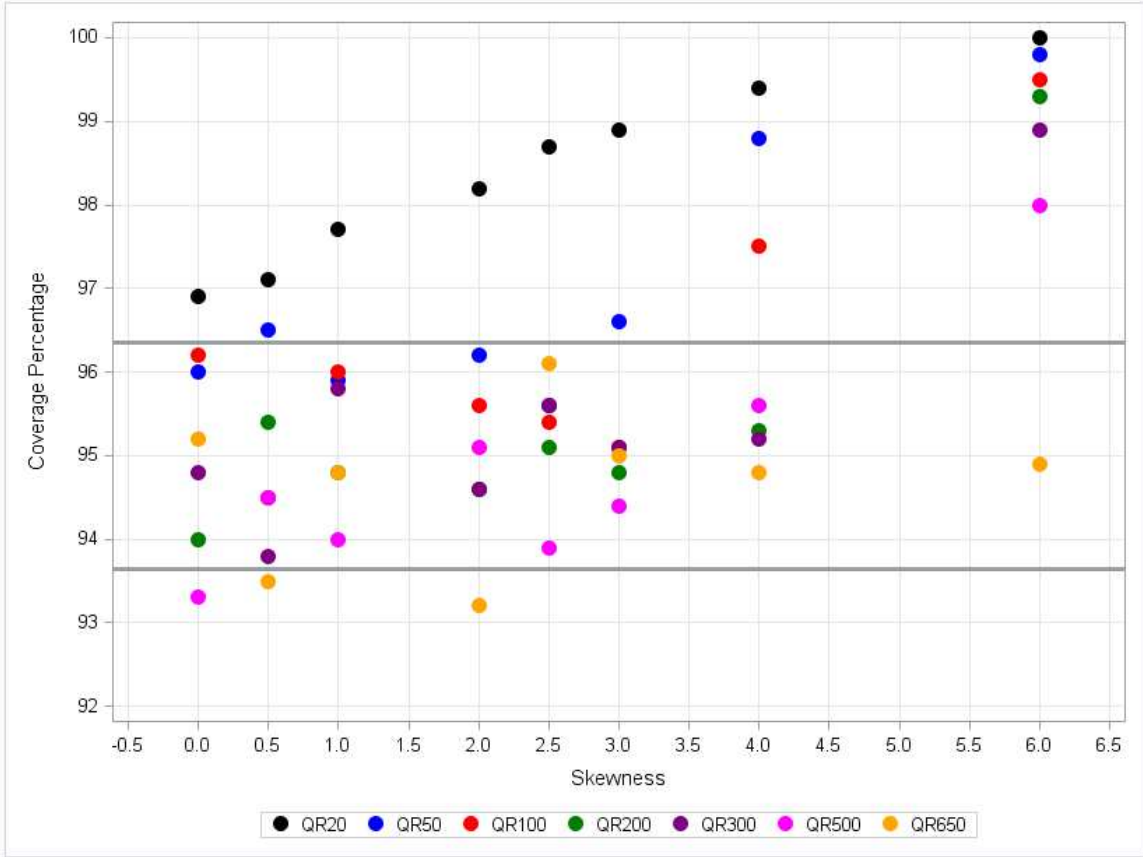


Figure 4.2: Coverage Percentage by Skewness and Sample Size for Median QR (Slope=0.4)

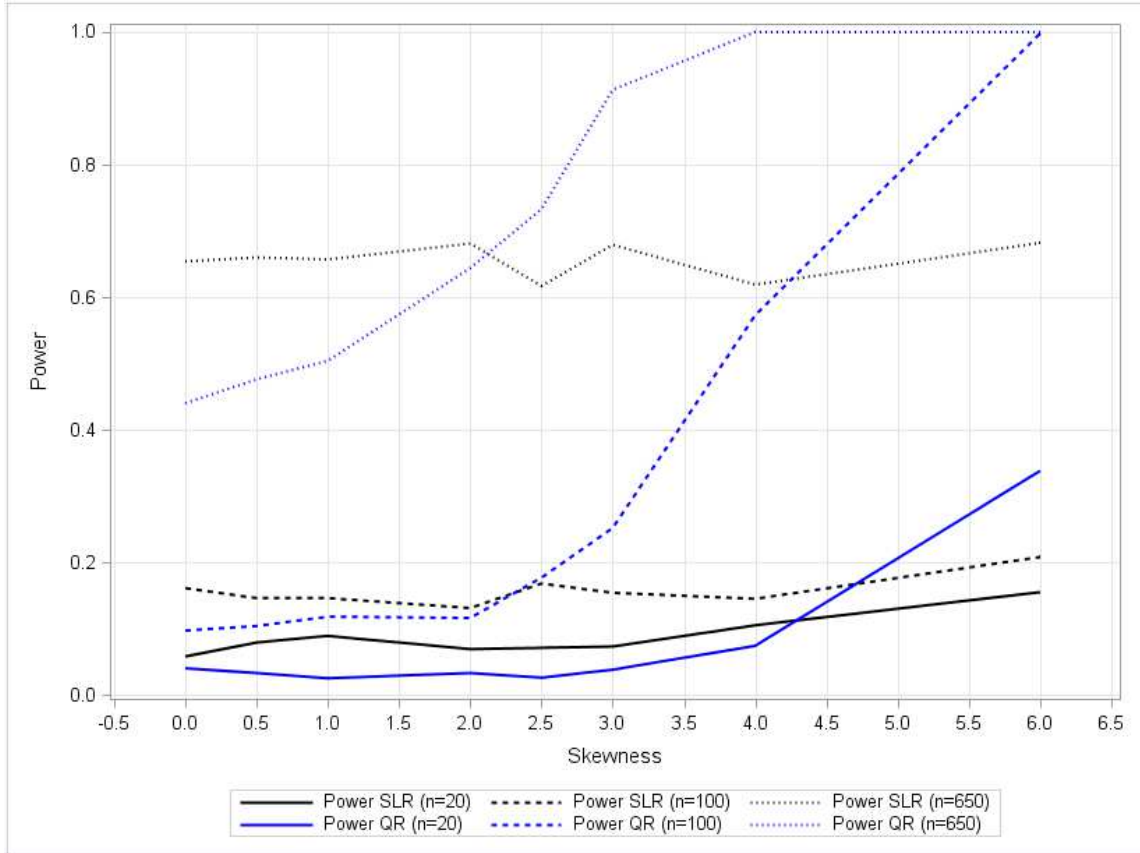


Figure 4.3: Power Comparison between SLR and Median QR(Slope=0.4) by Sample Size and Skewness of Error

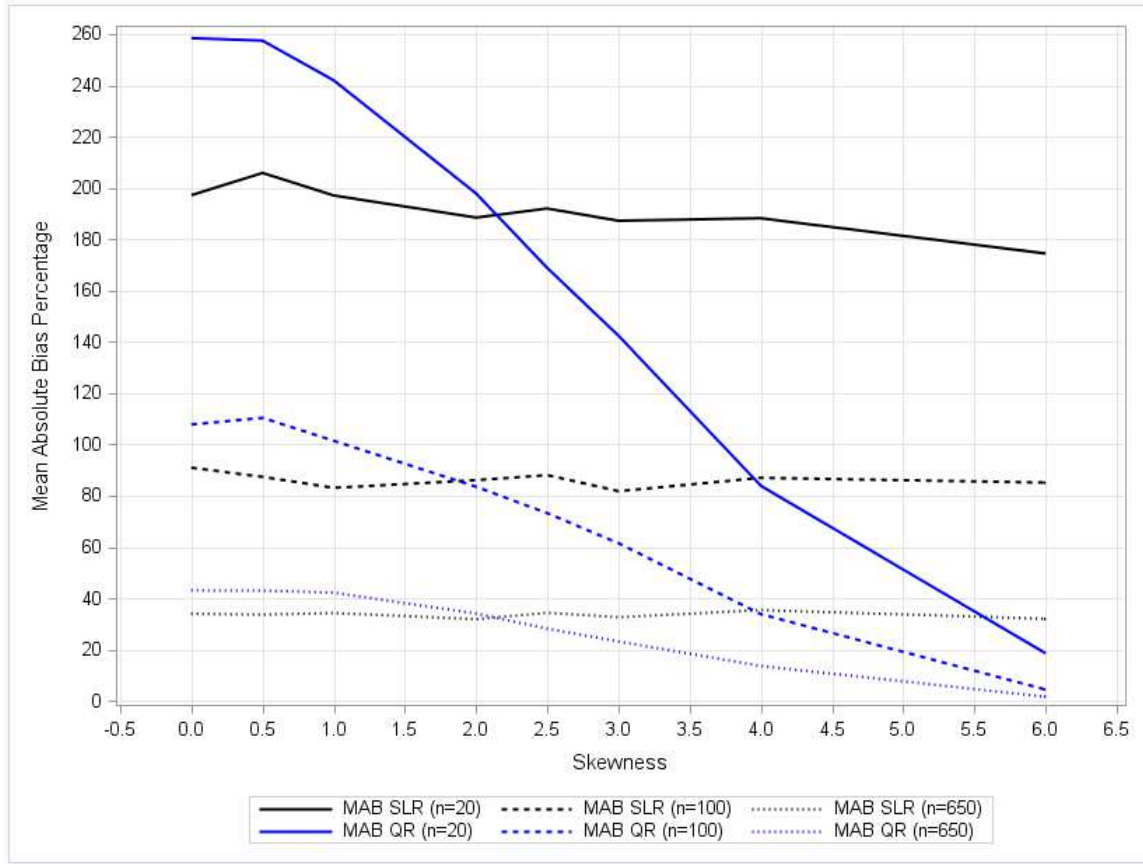


Figure 4.4: Comparison of Mean Absolute Bias(%) Between SLR and Median QR by Sample Size and Skewness of Error

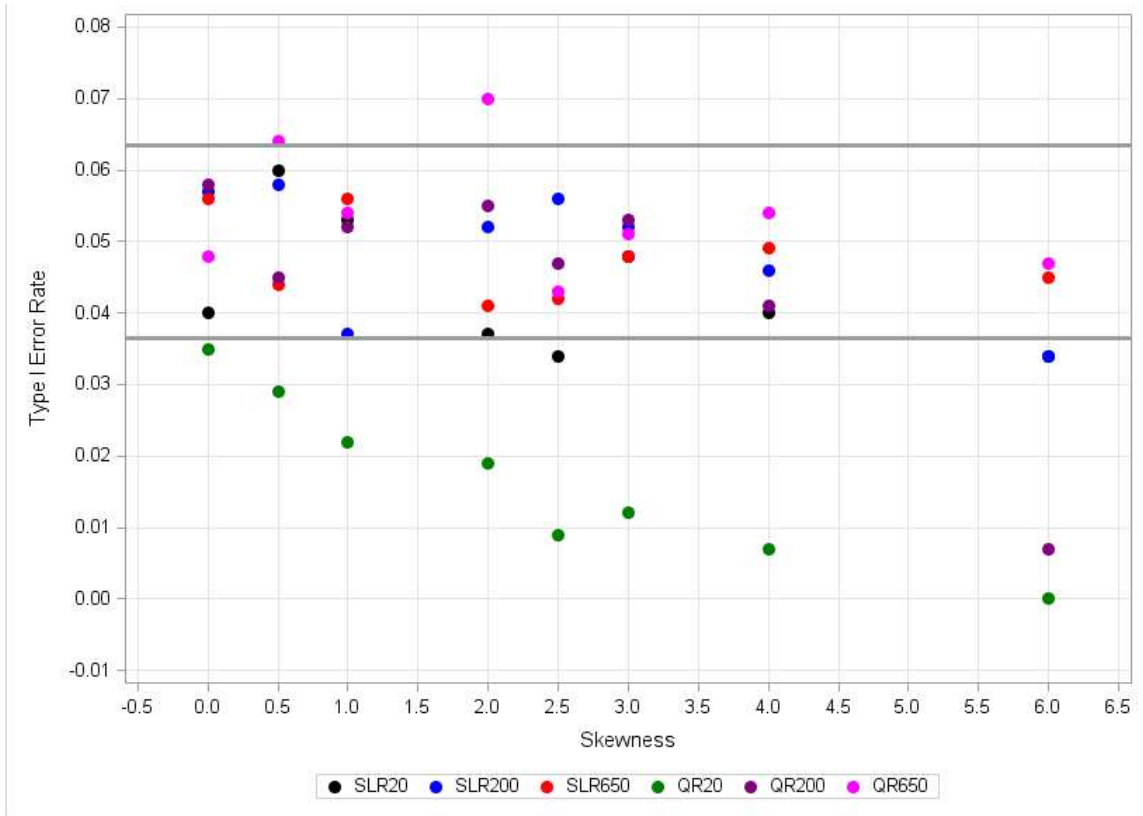


Figure 4.5: Type I Error Rate of SLR and Median QR by Sample Size and Skewness of Error

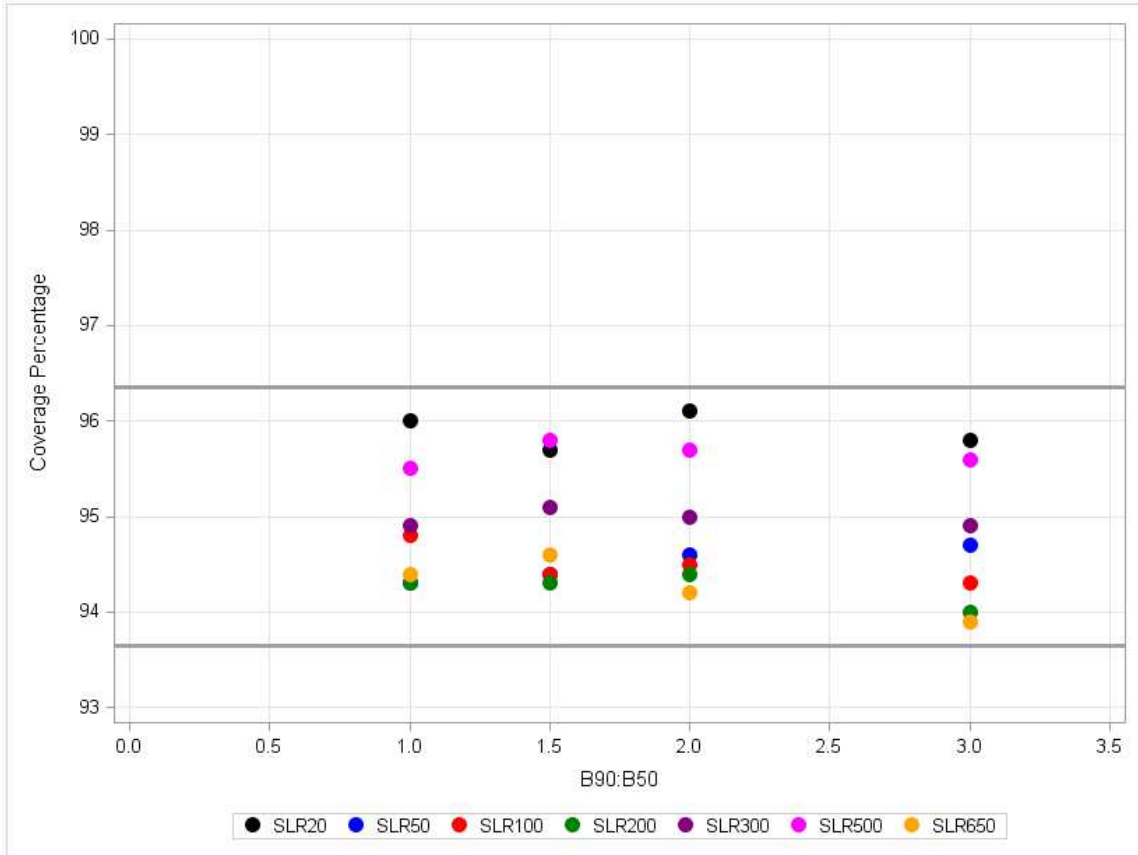


Figure 4.6: SLR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity

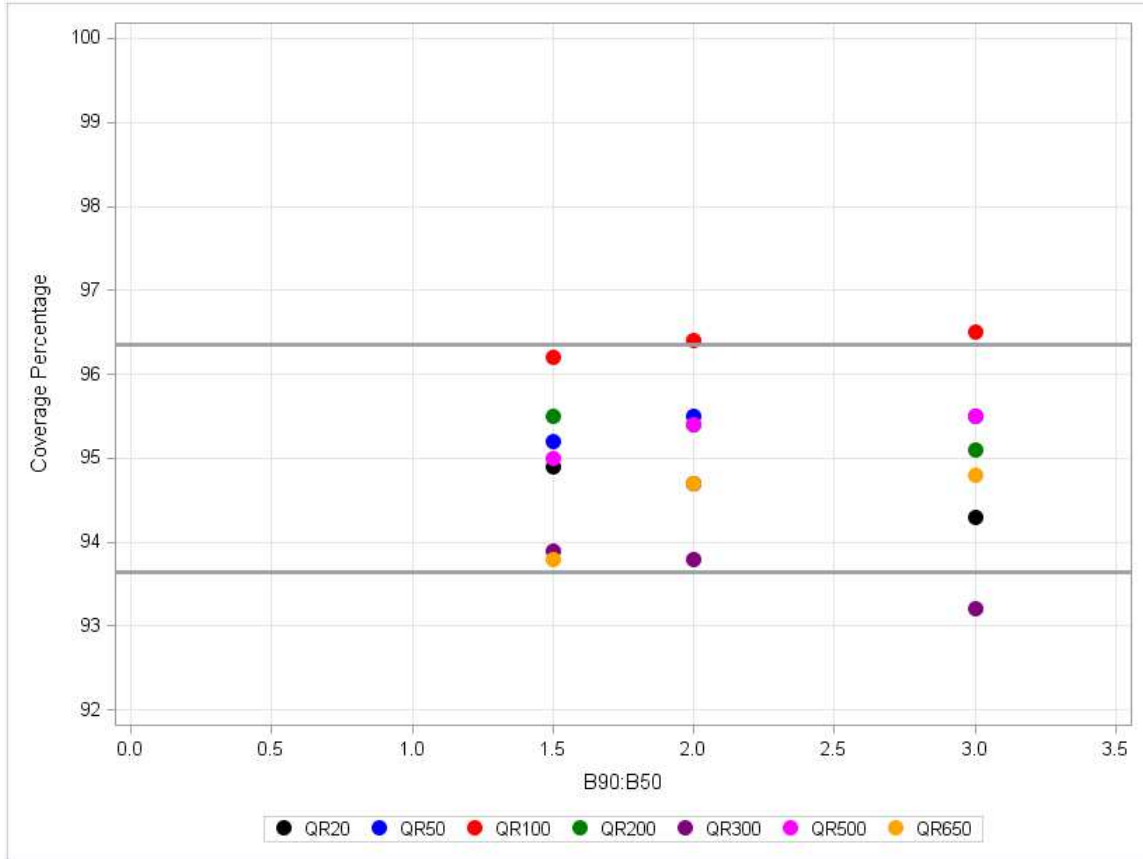


Figure 4.7: 10th Percentile QR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity

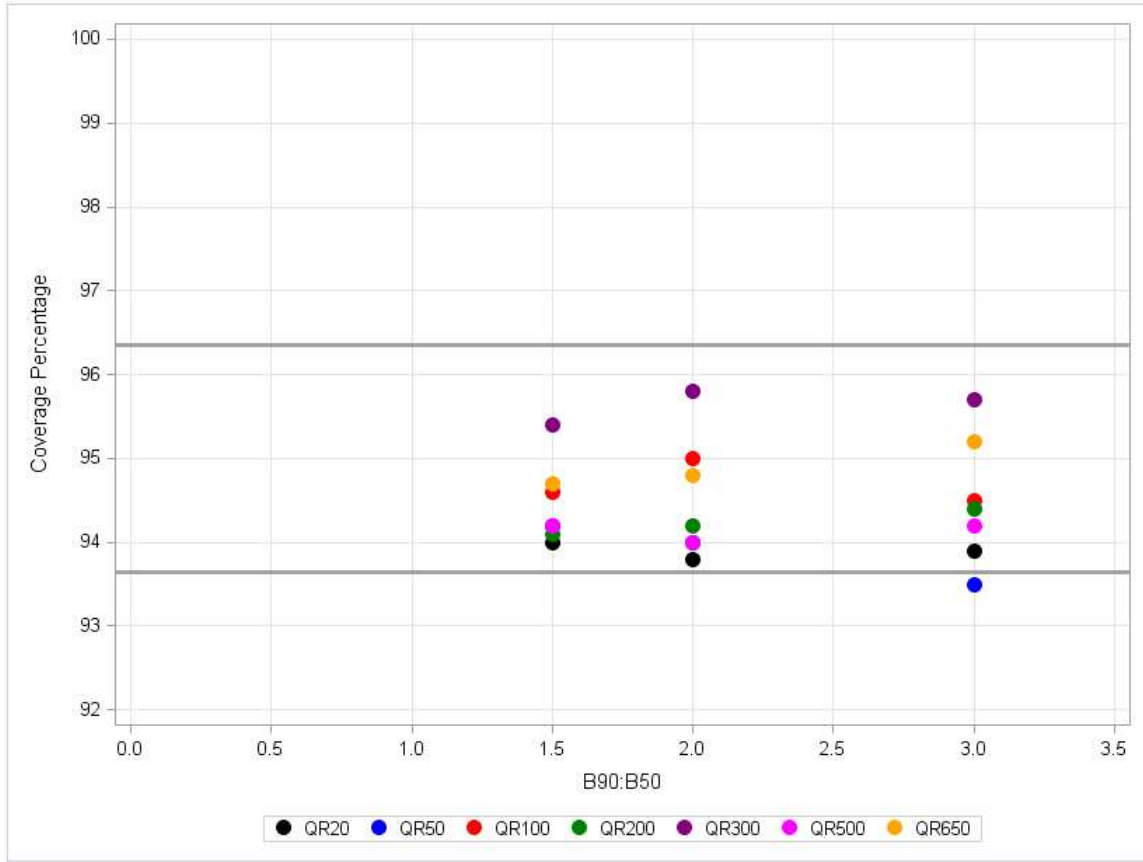


Figure 4.8: 90th Percentile QR Coverage Percentage (Slope=0.4) by Sample Size and Level of Heteroscedasticity

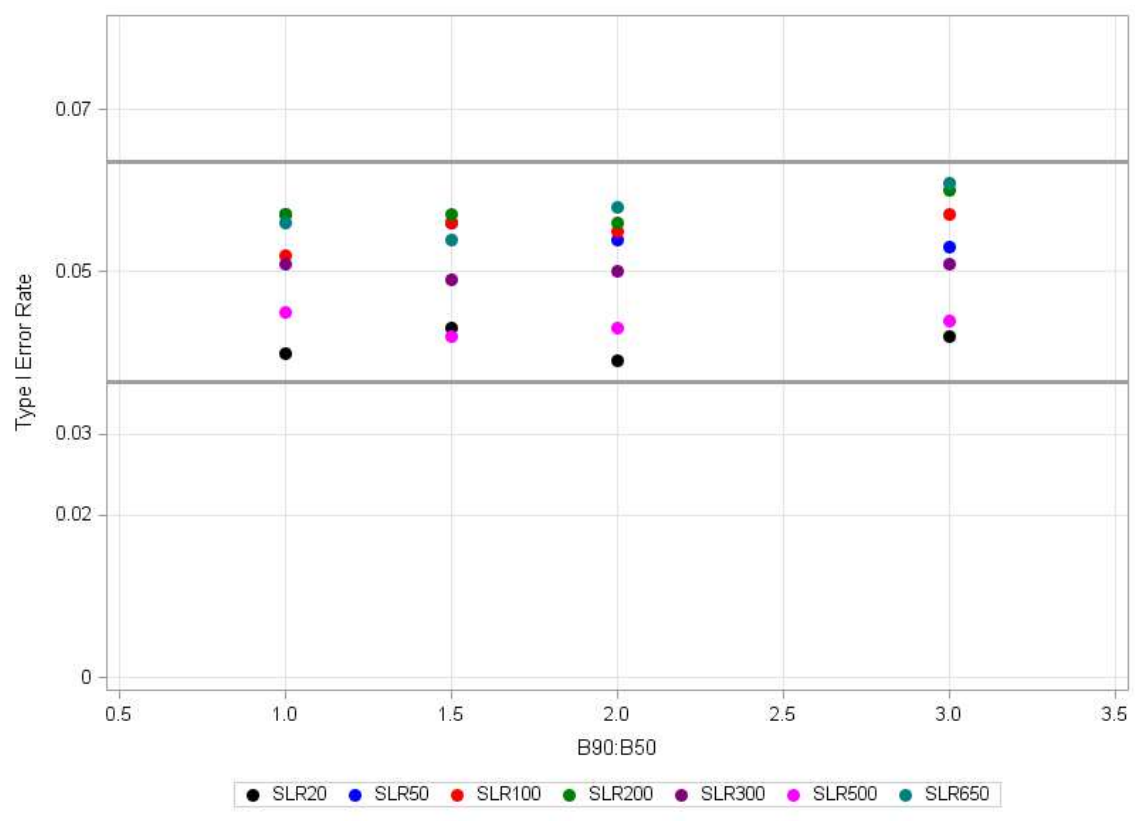


Figure 4.9: Type I Error Rate of SLR by Sample Size and Level of Heteroscedasticity

CHAPTER 5

CONCLUSIONS

Oftentimes, the research question itself can define whether quantile regression is desired over linear regression, as public health research is often more concerned with specific percentiles of the conditional distribution of an outcome variable rather than simply the mean. In cases where different methods of analysis could be used, it is therefore important to understand the consequences of performing statistical analyses when the assumptions are not met. Linear regression is a commonly used and easily understood method of analysis and, when the assumptions are met, it is able to provide a full description of the relationship between a predictor and outcome. When the assumptions of linear regression hold, specifically normality and homoscedasticity of the error distribution, the results provide an unbiased estimate of the relationship between the predictor and the outcome variable. In these cases, the slope representing the relationship between the predictor and outcome would be the same for any percentile. However, these assumptions are often not met in real-world scenarios, in which case inferences and the accuracy of the results can suffer. Research questions in the field of public health often involve variables that can be highly skewed or display variances that change as a function of the predictor. Quantile regression can provide an alternative to linear regression in times when the linear regression assumptions are not met. However, the possible penalty for relaxing the assumptions is a loss in power as compared to linear regression when these assumptions are met.

Our study suggests that, when the normality assumption is violated, specifically when the errors are skewed, there is a threshold at which quantile regression offers an improvement in statistical power over linear regression. After the skewness of the error distribution exceeds about 2, quantile regression not only becomes a more powerful method of analysis, but also it is able to more closely estimate the relationship between the predictor and the median of the outcome, as evidenced by the lower mean absolute bias percentage. When the sample size is small (e.g. less than 100), the loss of power when using quantile regression over linear regression is relatively small prior to this threshold (a difference in statistical power of 0.018 when skewness of error is 1 and sample size is 50), while there is a large potential gain in power when the skewness is high (a difference in statistical power of 0.187 when skewness of error is 4 and sample size is 50). However, our study suggests that when the sample size is even smaller ($n=20$) the skewness must be even greater before quantile regression surpasses linear regression in statistical power. This is likely due to the fact that such a small sample size leads to a lower likelihood of randomly selecting an extreme value. Regardless of sample size, our results suggest that, when the error distribution is not skewed, or even at very low levels of skewness, linear regression offers an improvement in statistical power. However, this improvement is smaller in magnitude than the improvement in power that quantile regression offers over linear regression if the error distribution is highly skewed.

The homogeneity of variance assumption of linear regression implies that the slope representing the relationship between the predictor and the outcome is approximately the same for any conditional percentile of the outcome. When this assumption is violated, however, there is a large difference between the estimates of the

slope for specific percentiles. In some cases, the direction of the association can change at different percentiles (e.g. a positive association within the 90th percentile, but a negative association within the 10th percentile), a phenomenon that linear regression is unable to detect. The field of public health is frequently more concerned with the extremes of the conditional distribution rather than the typical value, as the extremes are often where health concerns are most common or most pronounced. This simulation experiment suggests that using linear regression models in cases with heteroscedasticity can lead to large under or overestimations of the true association within these tails of the conditional distribution, perhaps even missing opposite effects at different tails of the distribution. This, in turn, could potentially lead to misleading conclusions and ineffective public health policies.

When both of these assumptions are violated, linear regression results can be heavily impacted, leading to misleading results. In this experiment using deviation from the mean CRP as the outcome, linear regression greatly overestimated the relationship between predictor and response for the majority of the conditional distributions, yet drastically underestimated the relationship within the 90th percentile. This suggests that, for data with both skewness and heteroscedasticity, linear regression is not an appropriate method of analysis as estimates are misleading and inference may not be valid. When considering multiple quantiles simultaneously, quantile regression was able to more accurately describe the relationship between the predictor and specific quantiles of the outcome while maintaining proper coverage percentage.

In conclusion, regression diagnostics such as residual plots and Q-Q plots can be used in combination with linear regression models to determine in practice whether the

normality or equal variance of errors assumptions have been violated. Our simulation study suggests that if the data are not skewed or even only slightly skewed (i.e. a skewness of less than about 2), linear regression has higher power than quantile regression. However, as skewness increases above that threshold, QR is a more advantageous method than SLR. When the assumption of homogeneity of variance in the error was violated in this study, the estimates from linear regression were unable to accurately describe the association at various percentiles. Linear regression estimates often over or underestimated the true effect within the extremes of the conditional distribution of the outcome. In cases where both the assumptions of normality and homoscedasticity of error are violated, however, the results from linear regression can not only incorrectly estimate the association between predictor and outcome variables at specific percentiles of the conditional distribution, but also the coverage percentage falls short, suggesting that inferences might not be valid. While transformations are a possible option to handle the violations of these assumptions, for data where there is larger deviation from these assumptions or multiple violations, quantile regression performs noticeably better than standard linear regression and offers the benefit of providing easily interpretable results.

REFERENCES

- Abramson, J. L., & Vaccarino, V. (2002). Relationship Between Physical Activity and Inflammation Among Apparently Healthy Middle-aged and Older US Adults. *Archives of Internal Medicine*, *162*(11), 1286.
doi:10.1001/archinte.162.11.1286
- Abrevaya, J. (2001). The effects of demographics and maternal behavior on the distribution of birth outcomes. *Economic Applications of Quantile Regression*, 247-257. doi:10.1007/978-3-662-11592-3_12
- Bind, M., Peters, A., Koutrakis, P., Coull, B., Vokonas, P., & Schwartz, J. (2016). Quantile Regression Analysis of the Distributional Effects of Air Pollution on Blood Pressure, Heart Rate Variability, Blood Lipids, and Biomarkers of Inflammation in Elderly American Men: The Normative Aging Study. *Environmental Health Perspectives*, *124*(8). doi:10.1289/ehp.1510044
- Hao, L. & Naiman, D. (2007) *Quantile Regression* (07-149). Sage Publications.
- Kleinbaum, D.G., Kupper, L. L., Muller, K.E., & Nizam, A. (2008). *Applied regression analysis and other multivariable methods* (4th ed.). Belmont, CA: Thomson, Brooks/Cole.
- Koenker, R., & Bassett, G., Jr. (Jan. 1978). Regression Quantiles. *Econometrica*, *46*(1), 33-50.

- Koenker, R., & Hallock, K. F. (2001). Quantile Regression. *Journal of Economic Perspectives*, 15(4), 143-156.
- Pagano, M., & Gauvreau, K. (2000). *Principles of Biostatistics* (2nd ed.). Cengage.
- Taylor, C., Rogers, G., Goodman, C., Baynes, R. D., Bothwell, T. H., Bezwoda, W. R., . . . Hattingh, J. (1987). Hematologic, iron-related, and acute-phase protein responses to sustained strenuous exercise. *Journal of Applied Physiology*, 62(2), 464-469. doi:10.1152/jappl.1987.62.2.464
- Tomaszewski, M. (2003). Strikingly Low Circulating CRP Concentrations in Ultramarathon Runners Independent of Markers of Adiposity: How Low Can You Go? *Arteriosclerosis, Thrombosis, and Vascular Biology*, 23(9), 1640-1644. doi:10.1161/01.atv.0000087036.75849.0b
- Wackerly, D. D., Mendenhall, W., & Scheaffer, R. L. (2008). *Mathematical statistics with applications* (7th ed.). Belmont, CA: Brooks-Cole.
- Weight, L. M., Alexander, D., & Jacobs, P. (1992). Strenuous Exercise: analogous to the acute-phase response? *Clinical Science*, 81(5), 677-683.
- Wong, N. D., Pio, J., Valencia, R., & Thakal, G. (2001). Distribution of C-Reactive Protein and Its Relation to Risk Factors and Coronary Heart Disease Risk Estimation in the National Health and Nutrition Examination Survey (NHANES) III. *Preventive Cardiology*, 4(3), 109-114. doi:10.1111/j.1520-037x.2001.00570.x