

## Project 2: Data Analysis on a Birth Weight Dataset

Trevor Carpenter, Christina De Cesaris, Michelle Tran  
STA 135: Multivariate Data Analysis  
Professor Maxime Pouokam  
University of California, Davis  
June 2, 2021

# I. Introduction

## II. Data Exploration

## III. Analysis

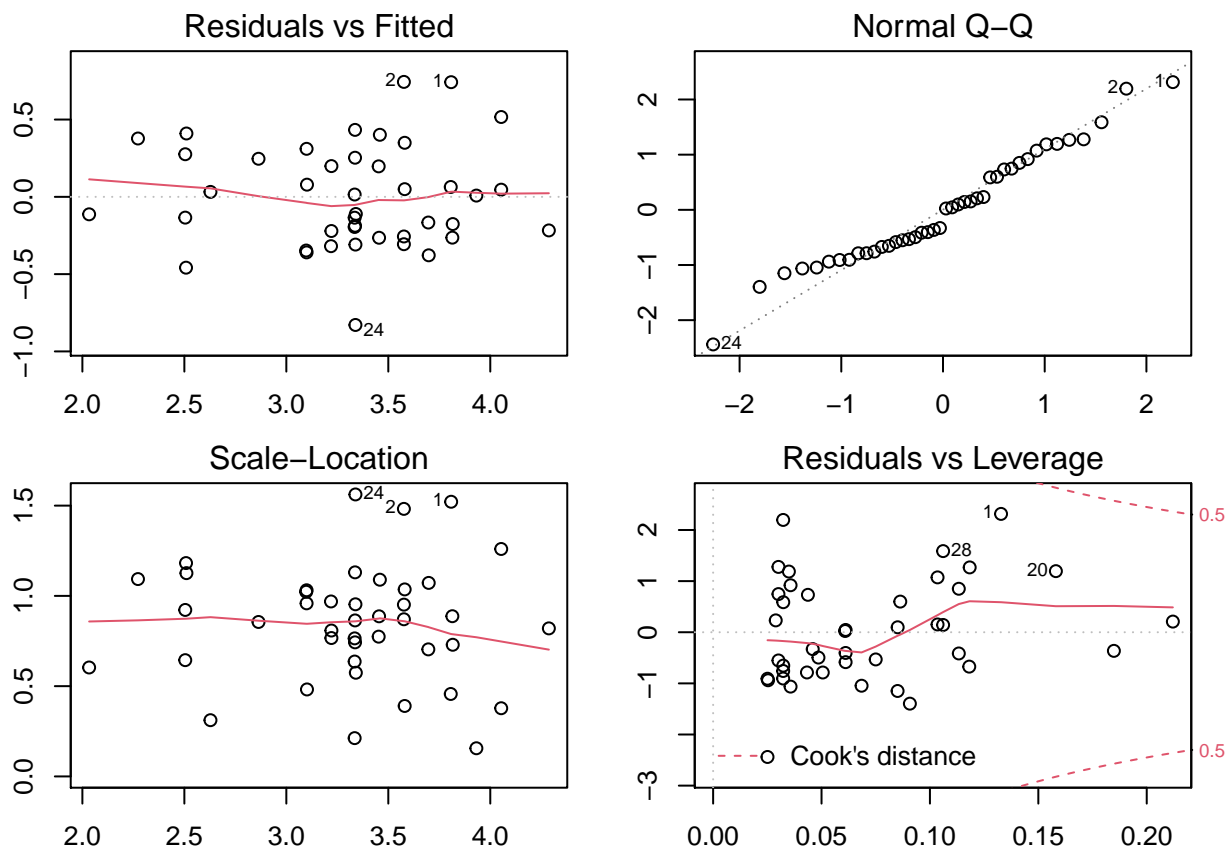
### i. Linear Model Fitting

Multiple linear regression attempts to model the relationship between two or more explanatory variables and a response variable by fitting a linear equation to observed data. In the given data, our response variable is the infant's birth weight and our explanatory variables are factors that supposedly effect birth weight. To get the reduced model, we looked at the full model's coefficients to determine if the predictor variable were significant. If it was significant, we can assume that changes in the predictor variable was significantly associated with changes in birth weight, and therefore its beta coefficient was not 0.

$$\text{Reduced Model: } Y = -5.45\beta_0 + .12\beta_2 + .118\beta_3$$

The reduced model involves the predictor variables: head circumference and gestation period.

After getting the reduced model, we must first check the model assumptions. If we do not check whether the model assumptions are true, it would lead to our model being imprecise due to possible outliers, non-constant variance, and nonlinearity of the outcome.



#### 1) Linearity of Data

For the Residual vs. Fitted plot, we can assume the linear assumption is reasonable because the line mostly lies on 0 and there are no patterns in the points.

#### 2) Normality of Residuals

Looking the QQ Plot above, although the model seems a little left skewed, because it is relatively

straight, we can assume the reduced model is normal. We can also check normality by also conducting a Shapiro-Wilks test. The result is that the p-value = 0.3, and therefore it fails to reject the normality null hypothesis. The reduced model's errors are normal.

3) Homogeneity of Residual Variance

For the Scale-Location plot, the red line is approximately horizontal and the residuals seems to be randomly scattered around the red line. This means that the spread of the residuals is roughly equal at all fitted values. We can assume that homoscedasticity is likely satisfied for the reduced model.

4) Outliers and High Leverage Points

Looking at the Residuals vs. Leverage plot, it highlights the top three extreme points (#1, #20, and #28). The most obvious outlier is point 20 at about 1.5 standard deviations above/below the mean. These outliers are, however, not influential because they are within the Cook's distance lines.

Overall, the reduced model appears to be adequate because it meets the model assumptions.

To test if our reduced model is correct in assuming that certain predictor variables are not significant, we can use the anova function to compare the full model against the reduced model to perform a partial F test.

$$H_0: \beta_1 = \beta_4 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = \beta_9 = \beta_{10} = \beta_{11} = \beta_{12} = 0$$

$$H_A: \beta_i \neq 0 \text{ for at least one } i$$

Using the partial F test, the F-statistics is  $F = 1.236$  and the p-value = 0.3106. This p-value is not significant at any  $\alpha$  value, and therefore we fail to reject the null hypothesis. We can conclude that the other predictor variables do not statistically contribute significantly to an infant's birth weight. The most statistically significant variables in determining an infant's birth weight are: head circumference and gestation period.

## ii. LDA & MANOVA

## IV. Conclusion

Linear Regression - two most significant variables in predicting birth weight were head circumference and gestation period and this is proven by our partial F test (rest are insignificant); reduced model is good bc we checked the 4 assumptions

## Appendix: R Script

```
knitr::opts_chunk$set(echo = F, warning = F, message = F)
smok = read.csv("data/Birthweight_reduced_kg_R.csv")
smok = smok[,2:14]
fmodel = lm(Birthweight ~., data = smok)

#stepwise
#step(fmodel)
#rmodel = lm(Birthweight ~ Headcirc + Gestation + smoker + mppwt,
#data = smok)
#summary(rmodel)

#Under the coefficient method for selecting reduced model, RSE is higher than using the step fcn
summary(fmodel)$coefficient
nmodel = lm(Birthweight ~ Headcirc + Gestation, data = smok)
summary(nmodel)
par(mfrow=c(2,2))
par(mar=c(2, 2, 2, 2))
plot(nmodel)

#normality test
ei <- nmodel$residuals
the.SWtest = shapiro.test(ei)

anova(fmodel,nmodel)
```