STAT 608, Spring 2022 - Assignment 3
SOLUTIONS

1. Question 1, Chapter 3, p. 103.

   (a) WHILE THIS MODEL DOES SEEM TO WORK WITH GREAT PRECISION, IT APPEARS THAT THE TRUE UNDERLYING RELATIONSHIP BETWEEN DISTANCE AND FARE MAY NOT BE A STRAIGHT LINE. IF, FOR EXAMPLE, THE UNDERLYING MODEL IS ACTUALLY A PARABOLA, OUR MODEL WILL SYSTEMATICALLY OVERESTIMATE OR UNDERESTIMATE THE PRICE OF THE AIRFARE, DEPENDING ON THE FLIGHT DISTANCE. THIS PROBLEM WILL BE EXACERBATED THE MORE WE ATTEMPT TO EXTRAPOLATE, IF WE DO.

   (b) AS ABOVE, IT IS LIKELY THAT WE NEED TO FIT SOMETHING OTHER THAN A STRAIGHT LINE RELATIONSHIP. IF WE USED A TRANSFORMATION OR PARABOLA, IT APPEARS WE MAY STILL FAIL TO FIT THE POINT WITH MAXIMUM AIRFARE; THAT POINT MAY NEED TO BE INVESTIGATED FOR OTHER EXPLANATIONS, AND PERHAPS AN INDICATOR VARIABLE CAN BE ADDED TO PREVENT IT FROM INFLUENCING THE FIT OF THE REST OF THE MODEL.

2. Explain in words why when we create confidence intervals and prediction intervals using a transformed response variable $Y$, we can't simply take the inverse transformation of the endpoints to get a confidence or prediction interval in the original units of $Y$.

   WHEN WE TRANSFORM THE DATA THE INTERVAL MAY CHANGE, AND WHEN WE TRY TO BACK-TRANSFORM THE INTERVAL, WE WILL NEED TO ADD OR MULTIPLY BY AN ADDITIONAL TERM IN ORDER FOR THE INTERVAL TO AGREE WITH THE ORIGINAL DATA. THIS IS KNOWN AS THE CORRECTION FACTOR AND IT CHANGES FOR A GIVEN TRANSFORMATION.

3. Recall the model with two indicator variables from question 3 of the previous homework. Calculate the hat matrix (use software if you like; it might be faster by hand). Explain what that projection matrix does and why it makes sense, as if to someone who has taken one semester of statistics.

   IN THIS DUMMY VARIABLE CASE, TAKE $m = 3$, $n = 4$, FOR EXAMPLE. THEN THE DESIGN MATRIX AND THE HAT MATRIX ARE

$$
\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, \ \mathbf{H} = \mathbf{X}\left(\mathbf{X'X}\right)^{-1}\mathbf{X'} = \begin{bmatrix} 1/3 & 1/3 & 1/3 & & & & \\ 1/3 & 1/3 & 1/3 & & & & \\ 1/3 & 1/3 & 1/3 & & & & \\ & & & 1/4 & 1/4 & 1/4 & 1/4 \\ & & & 1/4 & 1/4 & 1/4 & 1/4 \\ & & & 1/4 & 1/4 & 1/4 & 1/4 \\ & & & 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}
$$

   RECALL THAT $\hat{\mathbf{y}} = \mathbf{Hy}$. THE HAT MATRIX PROJECTS $\mathbf{y}$ ONTO ITS AVERAGE. THE DESIGN MATRIX TELLS US WHICH GROUP AN INDIVIDUAL IS IN: THE CALCIUM OR THE PLACEBO GROUP, SAY. SO THE HAT MATRIX WILL TAKE US FROM THE VECTOR $\mathbf{y}$, WHERE EVERY INDIVIDUAL HAS A DIFFERENT BLOOD PRESSURE, DOWN INTO THE VECTOR SPACE DEFINED BY THE COLUMNS OF $\mathbf{X}$, WHERE THERE ARE ONLY TWO GROUPS. SO THERE WILL ONLY

BE TWO BLOOD PRESSURES IN THAT SUB-SPACE: AN AVERAGE FOR THE CALCIUM GROUP AND AN AVERAGE FOR THE PLACEBO GROUP.

ALSO NOTE THAT THE TWO PARAMETERS OF INTEREST IN THIS MODEL ARE THE POPU-LATION MEANS FOR THE TWO GROUPS, AND OUR LINEAR MODEL HERE ESTIMATES THOSE TWO POPULATION MEANS WITH THEIR SAMPLE MEANS, AS WE MIGHT HAVE EXPECTED.

4. For the simple linear regression model in the case that our assumption is met that the errors are independent and identically distributed with variance $\sigma^2$:

   (a) Show that the formula for the vector of residuals $\hat{\mathbf{e}}$ can be expressed compactly using the hat matrix:
   $$(\mathbf{I} - \mathbf{H})\,\mathbf{y}$$
   IN THIS PROBLEM, THE MODEL $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ CAN BE WRITTEN AS
   $$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\mathbf{e}} \implies \hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{y} = \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y}.$$

   (b) Show that the covariance matrix of the residuals is therefore equal to
   $$(\mathbf{I} - \mathbf{H})\,\boldsymbol{\Sigma}\,(\mathbf{I} - \mathbf{H})',$$
   where $\boldsymbol{\Sigma}$ is the covariance matrix of the errors. Show that the covariance matrix of the residuals reduces to $(\mathbf{I} - \mathbf{H})\,\sigma^2$. (Please show that $\mathbf{H}$ is idempotent; that is, that $\mathbf{HH} = \mathbf{H}$.)
   IDEMPOTENT:
   $$\mathbf{HH} = \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}' = \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{I}\mathbf{X}' = \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}' = \mathbf{H}$$

   SYMMETRIC:
   $$\mathbf{H}' = \left(\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\right)' = \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}' = \mathbf{H}$$

   NOTE ALSO THAT
   $$(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) = \mathbf{I} - \mathbf{H} - \mathbf{H} + \mathbf{HH} = \mathbf{I} - \mathbf{H}$$
   THEREFORE
   $$\begin{aligned}\mathrm{Cov}\,(\hat{\mathbf{e}}|\mathbf{X}) &= \mathrm{Cov}\,((\mathbf{I} - \mathbf{H})\mathbf{y}|\mathbf{X}) = (\mathbf{I} - \mathbf{H})\mathrm{Cov}\,(\mathbf{y}|\mathbf{X})\,(\mathbf{I} - \mathbf{H})' \\ &= (\mathbf{I} - \mathbf{H})\sigma^2(\mathbf{I} - \mathbf{H})' = \sigma^2(\mathbf{I} - \mathbf{H})\end{aligned}$$

(c) Conclude that $\text{Cov}(\hat{e}_i, \hat{e}_j) = -h_{ij}\sigma^2$, $i \neq j$.

FROM (B), WHEN $i \neq j$, $\text{Cov}(\hat{e}_i, \hat{e}_j) = (0 - h_{ij})\sigma^2 = -h_{ij}\sigma^2$

5. For the simple linear regression model, show that the hat matrix $\mathbf{H}$ has the following properties:

   (a) $\mathbf{H}$ is symmetric.

$$\mathbf{H'} = \left( \mathbf{X}\left(\mathbf{X'X}\right)^{-1}\mathbf{X'} \right)' = \mathbf{X}\left(\mathbf{X'X}\right)^{-1}\mathbf{X'} = \mathbf{H}$$

   (b) $0 \leq h_{ii} \leq 1$, where $h_{ii}$ is the $i^{\text{th}}$ diagonal entry of the hat matrix. (**HINT**: First show that $h_{ii} \geq h_{ii}^2$ and note that $h_{ii} = \sum_j h_{ij}^2$.)

$$h_{ii} = \sum_j h_{ij}^2 = h_{ii}^2 + \sum_{i \neq j} h_{ij}^2 \geq h_{ii}^2 \implies h_{ii}(1 - h_{ii}) \geq 0 \implies 0 \leq h_{ii} \leq 1$$

   (c) The off-diagonals of the hat matrix are found by the formula

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX}.$$

FIRST REFERENCE TEXTBOOK SECTION 6.1.1, WHERE WE ARE TOLD THAT

$$\mathbf{X}\left(\mathbf{X'X}\right)^{-1}\mathbf{X'} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \frac{1}{SXX} \begin{bmatrix} \frac{1}{n}\sum_i x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} 1 & \cdots & 1 \\ x_1 & \cdots & x_n \end{bmatrix}$$

IF YOU WORK OUT THE ALGEBRA, YOU FIND THAT

$$\begin{aligned} h_{ij} &= \frac{1}{SXX}\left[ \frac{1}{n}\sum_i x_i^2 + x_i x_j - x_i \bar{x} - x_j \bar{x} \right] \\ &= \frac{1}{SXX}\left[ \frac{1}{n}\left( \sum_i x_i^2 - n\bar{x}^2 + n\bar{x}^2 \right) + x_i x_j - x_i \bar{x} - x_j \bar{x} \right] \\ &= \frac{1}{SXX}\left[ \frac{SXX}{n} + \bar{x}^2 + x_i x_j - x_i \bar{x} - x_j \bar{x} \right] \\ &= \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \end{aligned}$$

   (d) Finally, the text states that "There is a small amount of correlation present in standardized residuals, even if the errors are independent." Comment on when the covariances of the residuals are close to zero, for a fixed sample size. Why does it make sense that the covariances are close to zero in those situations?

WHEN $n \to \infty$, THE DENOMINATORS OF $h_{ij}$ WILL BE VERY LARGE AND $h_{ij} \to 0$. FOR A FIXED SAMPLE SIZE, IT MAKES SENSE THAT THE COVARIANCES ARE SMALL NUMBERS.

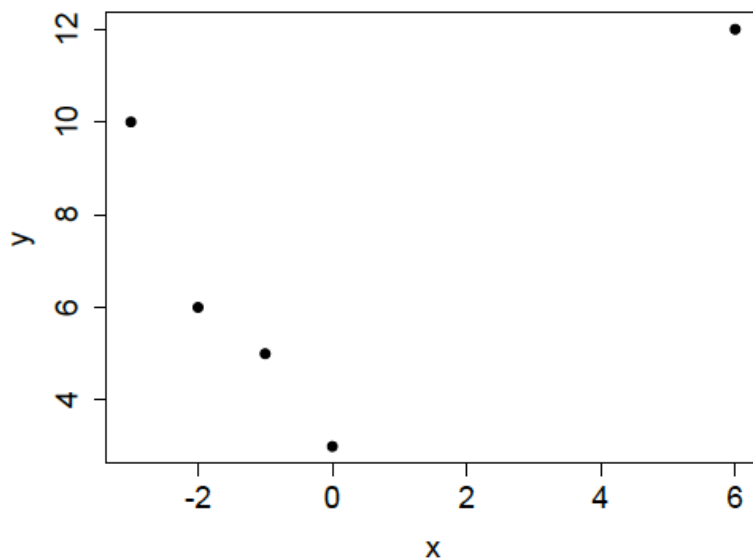Figure 1: Scatterplot of the data.

6. Under the simple linear regression model, $y_i = \beta_0 + \beta_1 x_i + e_i$, suppose the following data are collected and recorded: $x = [-3, -2, -1, 0, 6]$, $y = [10, 6, 5, 3, 12]$. Show your work, and do the following calculations by hand. (You may double-check with a computer.)

(a) First, create a quick sketch of the scatterplot of the data. Label any potential outliers or leverage points as such. Write out your design matrix.

   FIGURE 1 SHOWS THE SCATTERPLOT. THE DESIGN MATRIX IS

$$\mathbf{X} = \begin{bmatrix} 1 & -3 \\ 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 6 \end{bmatrix}$$

(b) Using $\hat{y} = 7.2 + 0.5x$, calculate the residuals $\hat{e}_i$.

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix} = 7.2 + 0.5 \begin{bmatrix} -3 \\ -2 \\ -1 \\ 0 \\ 6 \end{bmatrix} = \begin{bmatrix} 5.7 \\ 6.2 \\ 6.7 \\ 7.2 \\ 10.2 \end{bmatrix}$$

And

$$\hat{\mathbf{e}} = (\mathbf{y} - \hat{\mathbf{y}}) = \begin{bmatrix} 4.3 \\ -0.2 \\ -1.7 \\ -4.2 \\ 1.8 \end{bmatrix}$$

4

(c) Compute the leverage for each observation. Use the rule $h_{ii} > 4/n$ to identify potential leverage points. Are any points of high leverage "good" or "bad"?

$$\mathbf{X'X} = \begin{bmatrix} 5 & 0 \\ 0 & 50 \end{bmatrix}, \quad (\mathbf{X'X})^{-1} = \frac{1}{50} \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{H} = \mathbf{X}\left(\mathbf{X'X}\right)^{-1}\mathbf{X'} = \frac{1}{50} \begin{bmatrix} 19 & 16 & 13 & 10 & -8 \\ 16 & 14 & 13 & 10 & -2 \\ 13 & 12 & 11 & 10 & 4 \\ 10 & 10 & 10 & 10 & 10 \\ -8 & -2 & 4 & 10 & 46 \end{bmatrix}$$

USING THE RULE $h_{ii} > 4/n = 4/5 = 0.8$, WE CAN SEE $h_{55} = 46/50 > 0.8$, THE POINT $(6, \, 12)$ IS A POTENTIAL LEVERAGE POINT, AND IT IS "BAD" BECAUSE IT DOES NOT FALL NEAR THE LINE THAT WOULD FIT THE OTHER POINTS.

(d) Compute the variances of the residuals. (Assume the variance of the errors to simply be $\sigma^2$.)

THE VARIANCES OF THE RESIDUALS ARE

$$\mathrm{VAR}\,(\hat{e}_i) = \sigma^2(1 - h_{ii}) = \frac{\sigma^2}{50} \begin{bmatrix} 31 \\ 36 \\ 39 \\ 40 \\ 4 \end{bmatrix} = \sigma^2 \begin{bmatrix} 0.62 \\ 0.72 \\ 0.78 \\ 0.80 \\ 0.08 \end{bmatrix}$$

THE RESIDUAL $e_4$ HAS THE HIGHEST VARIANCE.

(e) Compute the standardized residuals $(s = 3.755)$. Comment no why this answer seems to conflict a bit with the answer to part (b) above.

THE STANDARDIZED RESIDUALS ARE

$$\frac{\hat{e}_i}{s\sqrt{1 - h_{ii}}} \begin{bmatrix} 1.454 \\ -0.063 \\ -0.513 \\ -1.251 \\ 1.695 \end{bmatrix}$$

POINT 5 HAS THE HIGHEST STANDARDIZED RESIDUAL IN ABSOLUTE VALUE, EVEN THOUGH POINT 1 HAS THE MOST EXTREME RESIDUAL. BECAUSE THE STANDARDIZED RESIDUALS ARE OBTAINED BY DIVIDING THE RESIDUALS BY THE VARIANCE OF THE RESIDUALS, POINTS WITH LOW VARIANCE WILL TEND TO HAVE HIGH STANDARDIZED RESIDUALS.

(f) Comment on why the point with the highest leverage in this dataset had the smallest variance.

SINCE THE VARIANCE IS GIVEN BY $\sigma^2(\mathbf{I} - \mathbf{H})$, POINTS WITH HIGH LEVERAGES WILL TEND TO HAVE LOW VARIANCES. THE BAD LEVERAGE POINT IS PULLING THE LINE TO BE VERY CLOSE TO IT, REDUCING THE VARIANCE OF THE RESIDUAL THERE.

7. When $Y$ has both mean and variance equal to $\mu$, we showed in the notes that the appropriate transformation of $Y$ for stabilizing the variance is the square root transformation. Now, suppose that $Y$ has mean equal to $\mu$ and variance equal to $\mu^2$. Show that the appropriate transformation of $Y$ for stabilizing variance is the log transformation. (Question 7, Chapter 3, page 112 of the textbook.)

SUPPOSE $f(Y) = \text{LOG}(Y)$, $\text{E}(Y) = \mu$, AND $\text{VAR}(Y) = \mu^2$. USING TAYLOR SERIES EXPANSION,

$$f(Y) = f(\mu) + f'(\mu)(Y - \mu) + \cdots$$

HENCE,

$$\text{VAR}(f(Y)) \approx \left(f'(\mu)\right)^2 \text{VAR}(Y) = \left(\frac{1}{\mu}\right)^2 \times \mu^2 = 1$$

BECAUSE THE VARIANCE OF $f(Y)$ IS CONSTANT IN $\mu$ (THAT IS, THE VARIANCE IS NO LONGER A FUNCTION OF THE MEAN), THE VARIANCE HAS BEEN STABILIZED AND THE LOG TRANSFORMATION IS APPROPRIATE.

8. Download the dataset called `company.csv` from Canvas. The dataset contains a systematic sample (every tenth company; we'll take these as randomly selected) for the Forbes 500 list. The variables of interest are `Sales` and `Assets` of the companies (both in millions of U.S. dollars). As with many financial datasets, many of these variables are skewed. Your job is to choose appropriate power transformations such that the relationship between `Assets` (response variable) and `Sales` (explanatory) are approximately linear.

   (a) Begin by creating a scatterplot of `Sales` and `Assets` and fit a simple linear regression line. What transformations does your scatterplot suggest? Create diagnostic plots for this model (Model 1). Discuss any weaknesses of this model.

   FIGURE 2 SHOWS THE SCATTERPLOT OF `Sales` AND `Assets`. IN FIGURE 3, THE Q-Q PLOT SHOWS THAT THE RESIDUALS ARE NOT NORMALLY DISTRIBUTED. THE FIRST AND THIRD PLOTS SHOW THAT MOST FITTED VALUES ARE ON THE LEFT-HAND SIDE OF THE GRAPH. THE SCALE-LOCATION PLOT SHOWS THAT OBSERVATIONS 16, 48, AND 54 ARE OUTLIERS; POINT 43 HAS A STANDARDIZED RESIDUAL THAT ALSO EXCEEDS OUR RULE OF THUMB OF ABSOLUTE STANDARDIZED RESIDUAL EXCEEDING 2. THE RESIDUAL VS. LEVERAGE PLOT SHOWS THAT OBSERVATION 40 IS A HIGH LEVERAGE POINT (ALL OF POINTS 33, 40, AND 43 HAVE LEVERAGE VALUES THAT EXCEED OUR RULE OF THUMB OF $4/n$). THE COOK'S DISTANCE PLOT IN FIGURE 4 SHOWS THAT OBSERVATION 16 IS THE POINT WITH THE WORST COOK'S DISTANCE; POINTS 40, 43, 48, AND 54 ARE ALSO ABOVE OUR CUTOFF FOR BAD OUTLIERS / HIGH INFLUENCE.

   (b) Choose an appropriate transformation for `Sales`. Explain how you made your choice. Include plots if applicable.

   I WOULD LIKE TO USE THE BOX-COX TRANSFORMATION TO SEE WHICH POWER IS THE MOST APPROPRIATE VALUE FOR DOING THE TRANSFORMATION. (WE KNOW THAT WHEN $\lambda = 0$, IT IS THE LOG TRANSFORMATION.)

   BY USING THE R CODE, WE FIND THAT $\lambda = -0.0675$ IS THE MOST APPROPRIATE POWER TRANSFORMATION, WHICH IS VERY CLOSE TO 0. ALSO NOTE THAT THE WALD
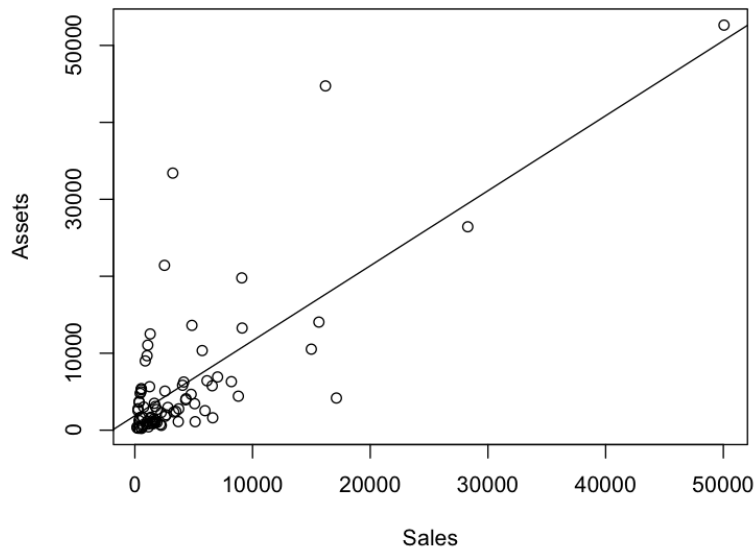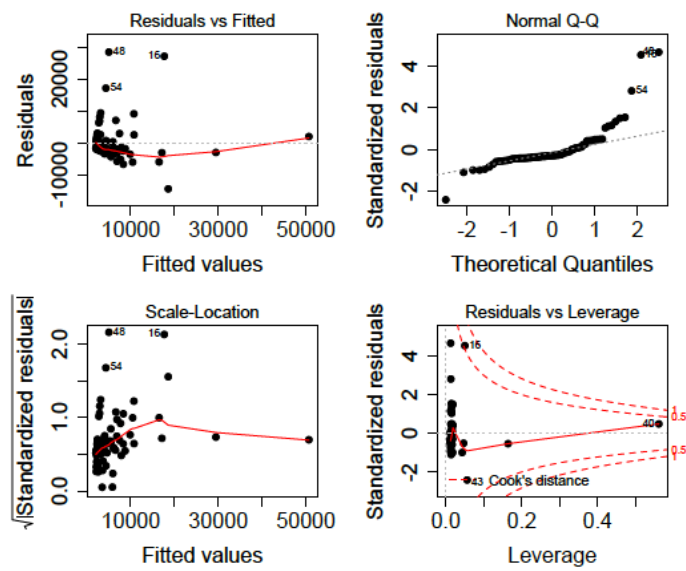
Figure 2: Scatterplot of the data
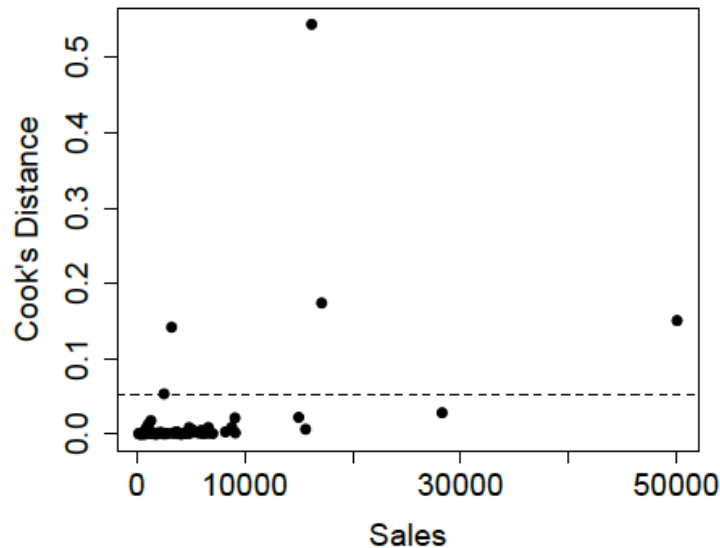


Figure 3: Diagnostic plot

7

Figure 4: Cook's distance plot

CONFIDENCE INTERVAL INCLUDES THE VALUE $0$, AND THE LIKELIHOOD RATIO TEST FOR TESTING THAT $\lambda = 0$ HAS A LARGE P-VALUE, SO WE FAIL TO REJECT THE HYPOTHESIS THAT A LOG TRANSFORMATION IS APPROPRIATE. THEREFORE, THE NATURAL LOG TRANSFORMATION FOR SALES IS AN APPROPRIATE TRANSFORMATION.

```
bcPower Transformation to Normality
      Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
Sales   -0.0675           0      -0.2329       0.0979


Likelihood ratio test that transformation parameter is equal to 0
 (log transformation)
                           LRT df     pval
LR test, lambda = (0) 0.6481734   1 0.42077


Likelihood ratio test that no transformation is needed
                          LRT df       pval
LR test, lambda = (1) 170.5833   1 < 2.22e-16
```

(c) Choose an appropriate transformation for `Assets`, and again explain how you made your choice. Because using an inverse response plot in this example is messy, you can just (1) fit a regression model of `Assets` vs. the transformed version of `Sales` that you chose in part (b), then (2) pass the fitted model into the `powerTransform` function. No plots required.

THIS TIME, I RUN A REGRESSION MODEL REGRESSING `Assets` ON $\log($`Sales`$)$. I USE THAT MODEL IN THE POWERTRANSFORM (BOX-COX) FUNCTION IN R SO THAT I TRANSFORM RESIDUALS TO BE AS CLOSE TO NORMAL AS POSSIBLE (RATHER THAN THE

8

$y$-VARIABLE). WE FIND THAT THE BEST TRANSFORMATION FOR $y$ IS NOW $-0.0166$, AND $0$ IS IN THE WALD CONFIDENCE INTERVAL; WE ALSO FAIL TO REJECT A NULL HYPOTHESIS THAT $\lambda = 0$ BECAUSE OF THE LARGE P-VALUE. THEREFORE THE NATURAL LOG TRANSFORMATION IS AN ACCEPTABLE TRANSFORMATION FOR $y$.

```
bcPower Transformation to Normality
   Est Power Rounded Pwr Wald Lwr Bnd Wald Upr Bnd
Y1   -0.0166        0      -0.1688        0.1357

Likelihood ratio test that transformation parameter is equal to 0
 (log transformation)
                             LRT df     pval
LR test, lambda = (0) 0.04564128   1 0.83083

Likelihood ratio test that no transformation is needed
                        LRT df       pval
LR test, lambda = (1) 159.6628   1 < 2.22e-16
```

(d) Call the model with both variables transformed Model 2. Create diagnostic plots for this model, and discuss any weaknesses of this model.

AFTER LOG TRANSFORMATIONS OF BOTH Sales AND Assets, THE NEW $x$ AND $y$ VARIABLES ARE LINEARLY RELATED, AS SHOWN IN FIGURE 5. THE DIAGNOSTIC PLOTS IMPROVE A LOT COMPARED WITH MODEL 1, AS SHOWN IN FIGURE 6. THE RESIDUAL PLOT IS LESS CONCENTRATED THAN MODEL 1, AND THE Q-Q PLOT SHOWS THE RESIDUALS ARE CLOSER TO A NORMAL DISTRIBUTION. IN MODEL 2, ALL ASSUMPTIONS ARE MET: LINEAR RELATIONSHIP, INDEPENDENT OBSERVATIONS, NORMALITY OF RESIDUALS AND CONSTANT VARIANCE.

(e) Compare Model 1 and Model 2. Which model is preferable?

AS SHOWN IN (D), MODEL 2 MET ALL THE ASSUMPTIONS. THEREFORE, I PREFER MODEL 2.

(f) Using the model $\log(\text{Assets}) = \beta_0 + \beta_1\log(\text{Sales})$, interpret the slope in the context of the problem.

THE SLOPE IS $0.587$, WHICH CORRESPONDS TO THE PERCENTAGE CHANGE IN ASSETS FOR EVERY PERCENTAGE CHANGE IN SALES. ACCORDING TO OUR MODEL, FOR EVERY $1\%$ INCREASE SALES, THERE IS APPROXIMATELY A $0.587\%$ INCREASE IN ASSETS, ON AVERAGE.

(g) Again using the model $\log(\text{Assets}) = \beta_0 + \beta_1\log(\text{Sales})$, find a 95% confidence interval for the average assets of a company with $6,571$ million in sales, as Hewlett-Packard did. Interpret your confidence interval in context.

TO TRANSFORM THIS INTERVAL INTO THE ORIGINAL UNITS, ADD ONE HALF OF THE MEAN SQUARE ERROR TO EACH ENDPOINT AND APPLY THE EXPONENTIAL FUNCTION, WE HAVE A $95\%$ CONFIDENCE INTERVAL OF $(6,870, 12,889)$. I AM $95\%$ CONFIDENT THAT COMPANIES WITH $6{,}571$ MILLION IN SALES WOULD HAVE ASSETS BETWEEN $6{,}870$ AND $12{,}889$ MILLION, ON AVERAGE.
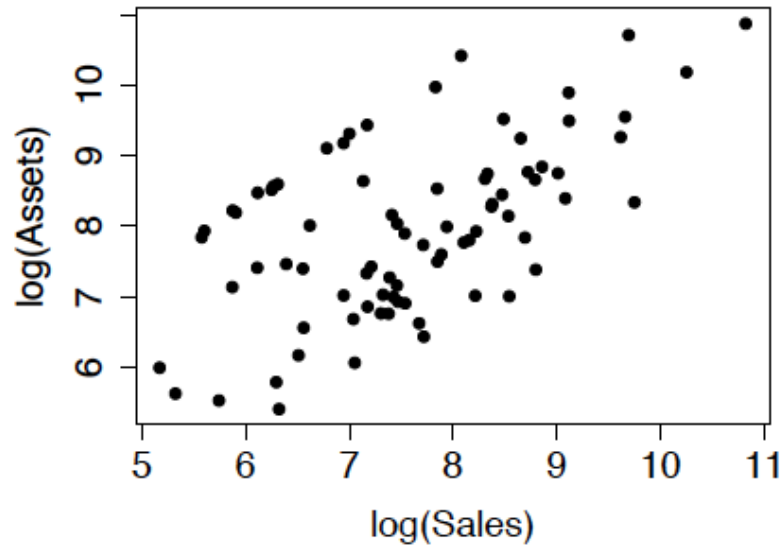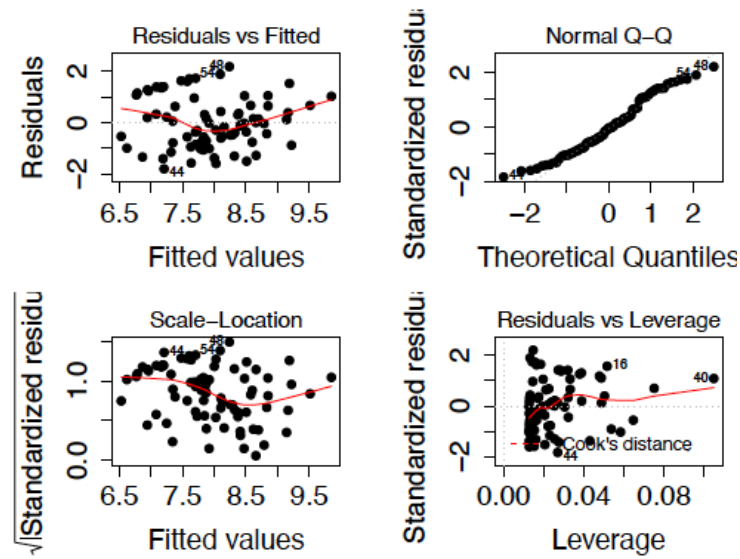
Figure 5: The relationship between `Sales` and `Assets` after transformation



Figure 6: Diagnostic plot after transformation