

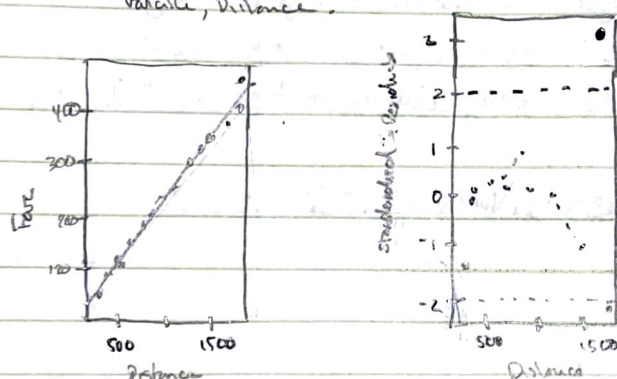
1.) Chp 3, Exercise 1 (A modern approach to Regression w/ R, pg. 103)

(1) The data file `airfare.txt` gives the one-way airfare (in US dollars) and distance (in miles) from city A to 14 other cities in the U.S. Interest centers on modeling airfare as a function of distance. The first model fit to the data was

$$\text{Fare} = \beta_0 + \beta_1 \text{Distance} + \epsilon$$

(a) Based on the output for the above model, a business analyst concluded the following:

The regression coefficient of the predictor variable, Distance, is highly statistically significant and the model explains 99.4% of the variability in the Y-variable, Fare. Thus the model is a highly effective model for both understanding the effects of Distance on Fare and for predicting future values given the value of the predictor variable, Distance.



Provide a detailed critique of these conclusions.

• First, I would note that the standardized residuals do not seem to randomly fluctuate around 0. In fact they seem to be a quadratic function of  $x$  (they seem to have a parabolic shape to them). This indicates to me that we haven't fit a correct model and should try improving the model by adding a quadratic term to our model.

• Second, we have two points (the two points where Distance > 1500) that can be classified as bad leverage points. I would first check to see if these two points are unusual or different in some way from the rest of the data. If we conclude that they are, I would remove them from our dataset and rerun our regression and see if the problem w/ the residuals (mentioned above) persists. If they are not different from the rest of the data, I would rerun the regression w/ the quadratic term included and see if those points are still bad leverage points.

(1) contd

(b) Does the ordinary straight line regression model seem to fit the data well?

IF not perfectly describes how the model can be improved.

- see discussion from part (a). I would try the following model

$$\text{Force} = \beta_0 + \beta_1 \text{Distance} + \beta_2 (\text{Distance}^2) + e$$

see stat 241 p. 11  
pg 56

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  $Y$ .

- Because in general  $E[g(Y)] \neq g(E(Y))$ . In other words, the Expected value of a function of  $Y$  is not, in general, equal to the function evaluated at the expected value of  $Y$ . Using a Taylor series expansion of  $g(y)$  about  $\mu_y$  we obtain:

$$X = g(Y) = g(\mu_y) + g'(\mu_y)(Y - \mu_y) + \frac{1}{2}g''(\mu_y)(Y - \mu_y)^2 + R$$
$$\mu_x = E[X] = E[g(Y)] = g(\mu_y) + g'(\mu_y)E[Y - \mu_y] + \frac{1}{2}g''(\mu_y)E[(Y - \mu_y)^2] + E[R]$$

$$\text{let } \bar{Y} = \mu_y$$

$$= g(\mu_y) + \frac{1}{2}g''(\mu_y)\text{Var}(Y) + E[R]$$



- 3) Recall the model w/ 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.

Recall Chp 3 Notes slide 13: hat matrix  $\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y = Hy$

Recall from Design 2, #3):

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix} \begin{matrix} \left. \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right\} m \\ \left. \begin{matrix} \vdots \\ \vdots \\ \vdots \end{matrix} \right\} n-m \end{matrix}$$

$$(X'X)^{-1} = \begin{bmatrix} \frac{1}{m} & 0 \\ 0 & \frac{1}{n-m} \end{bmatrix}$$

$$X(X'X)^{-1} = \begin{bmatrix} 1/m & 0 \\ 0 & 1/(n-m) \end{bmatrix} \begin{matrix} m \\ n-m \end{matrix}$$

Q: Is there a name for this type of matrix

$$X(X'X)^{-1}X' = \begin{bmatrix} 1/m & 0 & \dots & 0 \\ 0 & 1/(n-m) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1/(n-m) \end{bmatrix} \begin{matrix} m \\ n-m \\ \vdots \\ n-m \end{matrix}$$

note:  $Hy =$

$$\begin{bmatrix} 1/m & \dots & 1/m & 0 & \dots & 0 \\ \vdots & & \vdots & & & \vdots \\ 1/m & \dots & 1/m & 0 & \dots & 0 \\ \vdots & & \vdots & & & \vdots \\ 0 & \dots & 0 & 1/(n-m) & \dots & 1/(n-m) \\ \vdots & & \vdots & & & \vdots \\ 0 & \dots & 0 & 1/(n-m) & \dots & 1/(n-m) \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_m \\ \vdots \\ y_{n-m} \\ \vdots \\ y_n \end{bmatrix}$$

$\hat{y} = Hy$ , the hat matrix projects  $y$  onto the  $E[y|X]$ . This makes sense b/c for given values of our predictor variables, our best guess for  $y$  would be the average value of  $y$  when  $X = x$ .

Q: Is the above correct  $\forall H$ , or  $\forall X$  where  $X$  is a design matrix?  
I believe this is the case, as our regression model is defined as  $E[y|X]$ .

see last slide

4) For the simple linear regression model in the case that our assumption is met that the errors are iid w/ var  $\sigma^2$ :

(a) Show that the formula for the vector of the residuals  $\hat{e}$  can be expressed compactly using the notation:

$$(I - H) y$$

$$y = X\hat{\beta} + \hat{e} = X(X'X)^{-1}X'y + \hat{e} = Hy + \hat{e}$$

$$\Rightarrow \hat{e} = y - Hy = (I - H)y$$

see slide 37

(b) Show that the covariance matrix of the residuals is therefore equal to

$$(I - H) \Sigma (I - H)'$$

Where  $\Sigma$  is the covariance matrix of the errors. Show that the covariance matrix of the residuals reduces to  $(I - H) \sigma^2$ . Please show that  $H$  is idempotent.

That is, that  $HH = H$

$$HH = X(X'X)^{-1}X'X(X'X)^{-1}X' = XI(X'X)^{-1}X' = X(X'X)^{-1}X' = H$$

by (a)

we are to assume  $(I - H)$  is constant in this example? (i.e.  $\text{cov}(\hat{e})$  is only  $\text{cov}(\hat{e}|X)$ )

$$\text{cov}(\hat{e}) = \text{cov}((I - H)y) = (I - H) \text{cov}(y)(I - H)'$$

see slide 37

$$= (I - H) \text{var}(y)(I - H)' = (I - H) \Sigma (I - H)' = (I - H) \sigma^2 I (I - H)$$

$$= \sigma^2 (I - H)(I - H)' = \sigma^2 (I - H)(I' - H')$$

$$= \sigma^2 (II' - IH' - HI' + HH')$$

$$= \sigma^2 (I - H' - H + H) \quad (HH' = H \text{ by idempotent property})$$

$$\left[ \begin{array}{l} \text{* NOTE: } H' = (X(X'X)^{-1}X')' = (X'(X'X)^{-1})'X' = X((X'X)^{-1})'X' = X(X'X)^{-1}X' = H \\ \text{* recall: } (AB)' = B'A' \\ (A')' = A \end{array} \right]$$

$$\text{cov}(\hat{e}) = \sigma^2 (I - H)$$

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

$$\neq i, j: i \neq j, I_{ij} = 0 \therefore \text{cov}(\hat{e}_i, \hat{e}_j) = \sigma^2 (I_{ij} - h_{ij}) =$$

$$= \sigma^2 (0 - h_{ij}) = -h_{ij}\sigma^2$$

$$\boxed{\text{cov}(\hat{e}_i, \hat{e}_j) = -h_{ij}\sigma^2, i \neq j}$$



5) For the simple linear regression model, show that the hat matrix  $H$  has the following properties:

(a)  $H$  is symmetric ( $A' = A$ )

$$H' = (X(X'X)^{-1}X')' = (X')'(X(X'X)^{-1})' = X((X'X)^{-1})'X' = X(X'X)^{-1}X' = H.$$

NOTE: we are using the following properties of matrices in the above:

$$\bullet (AB)' = B'A'$$

$$\bullet (A^{-1})' = (A')^{-1}$$

(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  $h_{ii} = \sum_{j=1}^n h_{ij}^2$ )

$$\bullet h_{ii} = \sum_{j=1}^n h_{ij}^2$$

$$h_{ii} = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$$

$$[\text{note } h_{ij}^2 \geq 0 \Rightarrow \sum_{j \neq i} h_{ij}^2 \geq 0]$$

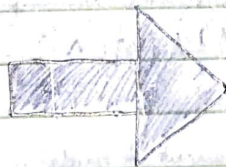
$$h_{ii} = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2 \geq h_{ii}^2$$

$$h_{ii} \geq h_{ii}^2$$

$$\bullet h_{ii} \geq h_{ii}^2 \Rightarrow h_{ii} - h_{ii}^2 > 0 \Leftrightarrow h_{ii}(1 - h_{ii}) \geq 0 \Rightarrow 0 \leq h_{ii} \leq 1$$

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

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_{k=1}^n (x_k - \bar{x})^2}$$



5.) (c) We know from previous work:  $(X'X)^{-1} = \frac{1}{SXX} \begin{bmatrix} \frac{SXX}{n} + \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix}$

$$(X'X)^{-1} X' = \frac{1}{SXX} \begin{bmatrix} \frac{SXX}{n} + \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}$$

$$= \frac{1}{SXX} \begin{bmatrix} \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_1 & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_n \\ -\bar{x} + x_1 & \dots & -\bar{x} + x_j & \dots & -\bar{x} + x_n \end{bmatrix}$$

$$X(X'X)^{-1}X' = \frac{1}{SXX} \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ x_i & x_i \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_1 & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_n \\ -\bar{x} + x_1 & \dots & -\bar{x} + x_j & \dots & -\bar{x} + x_n \end{bmatrix}$$

$$= \frac{1}{SXX} \begin{bmatrix} \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_1 - \bar{x}x_1 + x_1^2 & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j - \bar{x}x_1 + x_1x_j & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_n - \bar{x}x_1 + x_1x_n \\ \vdots & & \vdots & & \vdots \\ \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_1 - \bar{x}x_i + x_1x_i & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j - \bar{x}x_i + x_1x_j & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_n - \bar{x}x_i + x_1x_n \\ \vdots & & \vdots & & \vdots \\ \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_1 - \bar{x}x_n + x_1x_n & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j - \bar{x}x_n + x_1x_j & \dots & \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_n - \bar{x}x_n + x_n^2 \end{bmatrix}$$

$$h_{ij} = \frac{1}{SXX} \left( \frac{SXX}{n} + \bar{x}^2 - \bar{x}x_j - \bar{x}x_i + x_1x_j \right) = \frac{1}{n} + \frac{1}{SXX} (\bar{x}^2 - \bar{x}x_j - \bar{x}x_i + x_1x_j)$$

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

(d) Finally the text states that "There is a small amount of correlation present in the standard 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 these situations?

$$\text{cov}(e_i, e_j) = -h_{ij} \sigma^2 = -\frac{\sigma^2}{n} = \frac{\sigma^2 (x_i - \bar{x})(x_j - \bar{x})}{SXX}$$

$$= -\frac{\sigma^2}{n} = \frac{\sigma^2 (x_i - \bar{x})(x_j - \bar{x})}{\sum_{k=1}^n (x_k - \bar{x})^2}$$

• For a fixed  $n$  the  $\text{cov}(e_i, e_j)$  will be small when  $\text{var}(y|x=x) = \sigma^2$  is small.  
or if  $x_i$  is close to  $\bar{x}$  or if  $x_j$  is close to  $\bar{x}$ .

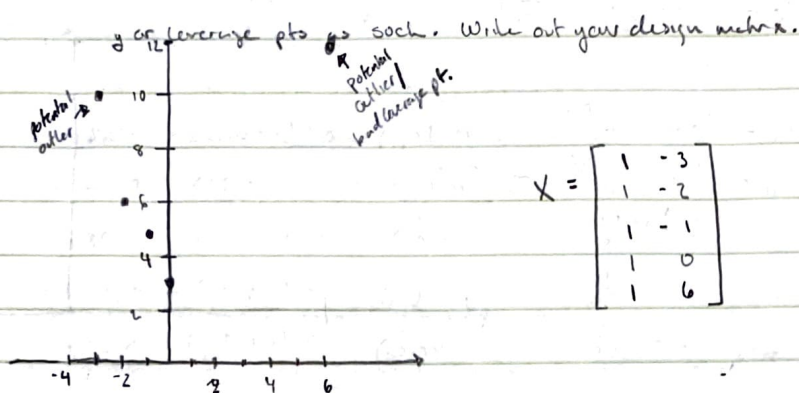
is this also  
var(e)  
yes. var(e|x=x) = var(y|x=x)  
2 mult.



$$\sum (x_i - \bar{x})^2$$

- 6.) Under the simple linear regression model,  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$  & the following data are recorded:  $x = [-3, -2, -1, 0, 6]$ ,  $y = [10, 6, 5, 3, 12]$ . Show your work and do the following calculations by hand.

(a) First create a quick sketch of the scatter plot of the data. Label any potential outliers



(b) Using  $\hat{y} = 7.2 + 0.5x$ , calculate  $\hat{\epsilon}$

$$\hat{\epsilon}' = [4.3, -0.2, -1.7, -4.2, 1.8]$$

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

Recall (Ch 3 Backs Slide 22):

- a bad leverage pt is a leverage pt whose standardized residual falls outside the interval  $[-2, 2]$
- a good leverage pt is a leverage pt whose standardized residual falls within the interval  $[-2, 2]$

leverage:  $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}$

leverage:  $h_{ii} = [0.38, 0.23, 0.22, 0.20, 0.42]$

cutoff for high leverage:  $\frac{4}{n} = 0.40$

- $(x, y) = (6, 12)$  is a high leverage pt.

standardized residual for  $(6, 12)$ :  $r_i = \frac{\epsilon_i}{\text{sd}(\epsilon_i)}$

$$s = \sqrt{\frac{1}{n-2} \sum \epsilon_i^2} = 3.755$$

$$r_i = \frac{1.8}{3.755 \sqrt{1-0.42}} = 1.7193 \approx 1.7$$

I know this is a bad question, but how do we get  $s$ ? is that not the correct way? - Never mind. we wanted  $s$  of the residuals, I was calculating  $s_y$  for  $\epsilon$ .

- yes, our 5th observation seems to be a "bad" leverage pt.

6) (contd)

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

$$\text{var}(\hat{e}_i) = \sigma^2 (1 - h_{ii})$$

Note: we know from (5) that for the simple linear regression model

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{XX}}$$

• For our case:  $h_{11} = \frac{1}{5} + \frac{(13 - 0)^2}{50}$

$$h_{22} = \frac{1}{5} + \frac{(12 - 0)^2}{50}$$

$$h_{33} = \frac{1}{5} + \frac{(11 - 0)^2}{50}$$

$$h_{44} = \frac{1}{5} + \frac{(10 - 0)^2}{50}$$

$$h_{55} = \frac{1}{5} + \frac{(6 - 0)^2}{50}$$

$$\text{var}(\hat{e}_1) = \sigma^2 (1 - 0.38) = 0.62 \sigma^2$$

$$\text{var}(\hat{e}_2) = \sigma^2 (1 - 0.28) = 0.72 \sigma^2$$

$$\text{var}(\hat{e}_3) = \sigma^2 (1 - 0.22) = 0.78 \sigma^2$$

$$\text{var}(\hat{e}_4) = \sigma^2 (1 - 0.20) = 0.80 \sigma^2$$

$$\text{var}(\hat{e}_5) = \sigma^2 (1 - 0.12) = 0.88 \sigma^2$$

(e) compute the standardized residuals:

(recall H.D. pg 20):  $r_i = \frac{\hat{e}_i}{\sqrt{1 - h_{ii}}}$

$$r' = [1.454, -0.063, -0.513, -1.251, 1.695]$$

• This conflicts a bit w/ (b) b/c pt 1 has the highest residual (calculated in (b)), but point 5 has the largest standardized residual. This occurs b/c pt 5 has a lower variance than pt 1.

(F) comment on why the pt w/ the largest leverage in this dataset had the smallest variance.

• Variance is given by  $\sigma^2 (1 - h_{ii})$  where  $0 \leq h_{ii} \leq 1$  is our leverage.

Thus, pts w/ higher leverage will tend to have lower variance.



- 7) When  $Y$  has both mean & variance equal to  $\mu$ , we showed in the notes that the appropriate transformation of  $Y$  for stabilizing the variance is the square root of the transformation. Now if  $Y$  has mean  $\mu$  and variance  $\mu^2$ , show the appropriate transformation of  $Y$  for stabilizing variance is the log transformation. (Q7 Chp 3, pg 112 in textbook).

[\* see chp 3 notes pgs 45, 46]

• let  $F(Y) = \ln(Y)$ ,  $E[Y] = \mu$ ,  $\text{var}[Y] = \mu^2$

•  $\text{var}(F(Y)) = (F'(E[Y]))^2 \text{var}(Y)$

•  $F'(\mu) = \frac{1}{\mu} \Rightarrow (F'(\mu))^2 = \frac{1}{\mu^2}$

$\Rightarrow \text{var}(F(Y)) = \frac{1}{\mu^2} \cdot \mu^2 = 1$  + constant wrt  $\mu$

- 8.) Download the dataset called company.csv from canvas. The data contains a systematic sample (every 10<sup>th</sup> company) we'll take time as randomly selected for the Forbes 500 list. The variables of interest are Sales & Assets of the companies. As w/ many financial datasets, many of these variables are skewed. Your job is to choose appropriate power transformations p.t. the relationship between Assets (response var) & Sales (EV) are approximately linear.

(a) The scatter plot seems to suggest a log transformation might be appropriate

- First weakness of the simple linear regression model we fit is that the residuals aren't normally distributed. This can be seen clearly in the Q-Q plot of the standardized residuals.
- We also have a few outliers in the data, specifically the 116<sup>th</sup> & 148<sup>th</sup> observations have standardized residuals  $> 2$ .
- Additionally we have 5 observations that seem to be highly influential (14, 40, 43, 98, 54) when looking at Sales vs Cooks Distances. All these points are above the cutoff  $> 1/2$  to be considered influential.
- Finally, we see that we have non-constant error variance.

see stats 411 notes 97-98 for box cox code

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

- Taking log(sales) seems to be an appropriate transformation. If run a box cox transformation and find the appropriate power to be  $-0.068$ . However a 95% CI for the appropriate power was  $(-0.235, 0.095)$ . Since the CI contained 0, used the log transform. Furthermore, a Shapiro-Wilk's test on the log transformed sales data gives a p-value  $> 0.05$  indicating we cannot conclude the log transformed data does not have a normal distribution.

- 8.) (c) Choose an appropriate transformation for assets, and again explain how you made your choice. B/c using an inverse response plot in this example is messy, you can (1) just fit a regression model of Assets vs. the transformed model of sales, then (2) pass the fitted model into the `powerTransform` function. No plots required.

I first did the same thing for Assets that I did for sales, and again found that the log transform worked. CI for  $\lambda_{\text{Assets}}$  (-0.203, 0.115) w/  $\hat{\lambda}_{\text{Assets}} = -0.043$

Fitting a regression model of Assets vs  $\log(\text{sales})$ , passing through the fitted model into `powerTransform` we get  $\hat{\lambda} = -0.01658166$ . A 95% CI for  $\lambda_{\text{Assets}}$  is given by (-0.1658, 0.1333). Thus, the log transformation seems appropriate (b/c our CI contains 0).

is the process normalizing the residuals? and the above is normalizing  $y$  independent of  $x$ ? both work

b/c when normalizing residuals, we get the st we want, when normalizing  $y$  independent of  $x$ ,  $x, y$  are apart of the same location scale family so we get almost the same. (To see a way to see if we should normalize just one or the other, or both simultaneously)

(d) Call the model w/ both variables transformed model 2. Create diagnostic plots for this model; discuss any weaknesses of this model?

After the log-log transformation, sales & assets seem to be linearly related.

The diagnostic plots are not perfect, but have definitely improved over the original model w/ the untransformed variables. The residuals vs fitted plot is less curved - seem to have constant variance, the Normal Q-Q plot shows our residuals might have lighter tails than if they were normally distributed, but they are closer to normal than they were before we transformed the variables.

When looking at the plot of sales vs coasts distance, we still seem to have 2 influential pts (16, 40).

Looking at the plot of fitted values vs (standardized residuals)<sup>1/2</sup> we seem to have constant error variance.

Overall model 2 seems to be an improvement over model 1 w/ it meeting the assumptions of a regression model.

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

As discussed in (d) model 2 better meets the assumptions of a linear regression model. Thus, I prefer model 2.

see stat 641 H.O. 9.19 42 for CI

ask how to get CI for  $\lambda$  from `powerTransform` function in R

use `summary(powerTransform)`



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

• Model 2 gives us  $\hat{\beta}_1 = 0.5870$ . In the context of this model, that means <sup>on the Forbes 500 list.</sup>  
For a 1 percentage point increase in the sales of a company, we expect, on average, to see a corresponding 0.5870 percentage point increase in the assets of a company in the Forbes 500 list.

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

NOTE: We transformed the data so we must apply a back transformation.  
• Ask about calculation of MSE. relative to Residual standard error in R.

I am 95% confident that a company on the Forbes 500 list w/ 6,571 million in sales would on average have assets between 6,870 and 12,889 million.

Q. Do we just back transform the mean, the do  $\pm 1.96 \times \text{SE}$  or do we back transform the whole pt? See Chp 2 pg 65. Chp 3 pg 80.

