

# Homework 5

Aiden Kenny

STAT GR5205: Linear Regression Models

Columbia University

November 25, 2020

## Question 1

*Collaborators:* None

- (a) Let  $Y$  be the number of nurses in the hospital and let  $X$  be the available faculty and services. The left and middle panels of Figure 1 show the histograms of  $Y$  and  $X$ , respectively. We see that  $Y$  is skewed right, while  $X$  appears to be normally-distributed. In addition, the scatterplot of  $Y$  vs.  $X$ , which is in the third panel of Figure 1, shows that there is a nonlinear relationship between  $Y$  and  $X$ . All of these indicate that  $Y$  is suitable for a data transformation. Specifically, we would like to perform a power transformation on  $Y$  in order to make it closer to a normal distribution.



Figure 1: Histograms of  $Y$  and  $X$  and a scatterplot of  $Y$  vs.  $X$ .

- (b) The power transformation function and its scaled counterpart are defined as

$$f_{\lambda}(Y) = \begin{cases} Y^{\lambda} & \text{if } \lambda \neq 0 \\ \log Y & \text{if } \lambda = 0. \end{cases} \quad \text{and} \quad g_{\lambda}(Y) = \begin{cases} (Y^{\lambda} - 1)/\lambda & \text{if } \lambda \neq 0 \\ \log Y & \text{if } \lambda = 0. \end{cases}$$

When transforming  $Y$ , we first use the scaled power transform to fit the model  $g_{\lambda}(Y) = \beta_0 + \beta_1 X + \epsilon$  in order to determine an optimal value of  $\lambda$  via maximum likelihood estimation. We then use  $f_{\lambda}$  to make  $Y$  closer to a normal distribution and perform any subsequent inferences. In vector form, our model is  $\mathbf{g}_{\lambda}(\mathbf{y}) = \beta_0 \mathbf{1} + \beta_1 \mathbf{x} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$  and  $\mathbf{g}_{\lambda}(\mathbf{y})$  is the transformed response (so  $[\mathbf{g}_{\lambda}(\mathbf{y})]_i = g_{\lambda}(y_i)$ ) for some unknown  $\lambda$ . It can be shown (see Appendix A) that the log-likelihood function can be expressed as a function of only  $\lambda$ ,

$$m(\lambda) := -\frac{n}{2} \log \left( \frac{2\pi e}{n} \right) - \frac{n}{2} \log \left( \mathbf{g}_{\lambda}^T(\mathbf{y})(\mathbf{I} - \mathbf{H})\mathbf{g}_{\lambda}(\mathbf{y}) \right) + (\lambda - 1) \sum_{i=1}^n \log(y_i),$$

where  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is the hat matrix.  $m(\lambda)$  has been plotted in the left panel of Figure 2, where we can see that the MLE is maximized at  $\lambda = 0.085$ .

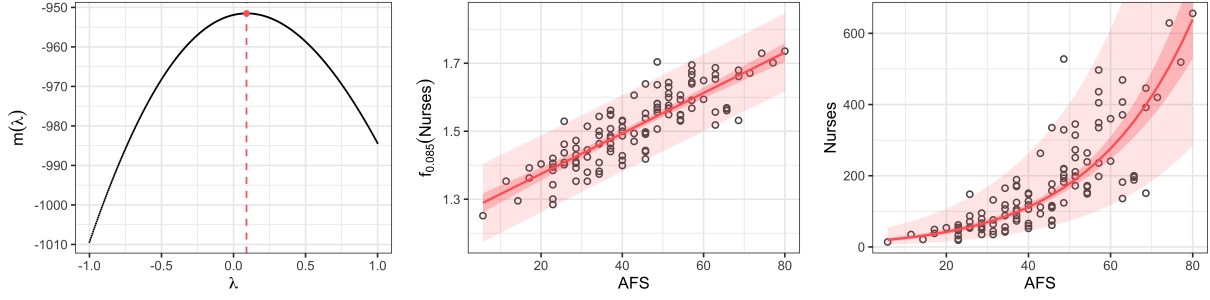


Figure 2: Relevant plots for the power transformation of  $Y$ .

- (c) Now that we have our optimal  $\lambda$ , we will fit the model  $f_{0.085}(Y) = \beta_0 + \beta_1 X + \epsilon$ . We have  $\hat{\beta}_0 = 1.255$  and  $\hat{\beta}_1 = 0.005953$ . A plot of this model, along with a 95% confidence interval, has been printed in the middle panel of Figure 2. This means that when  $X$  increases by 1 unit,  $f_{0.085}(Y)$  is expected to increase by 0.005953 units. We can also make inferences about  $Y$  itself; we have  $Y = f_{0.085}^{-1}(\hat{\beta}_0 + \hat{\beta}_1 X) = (1.255 + 0.005953X)^{1/0.085}$ . A plot of this model has been printed in the right panel of Figure 2. We must keep in mind that this model is non-linear, so it's expected rate of change will depend on the value of  $X$ .
- (d) Both diagnostic plots can be found in Figure 3, and the data seems to confirm the model assumptions very well.

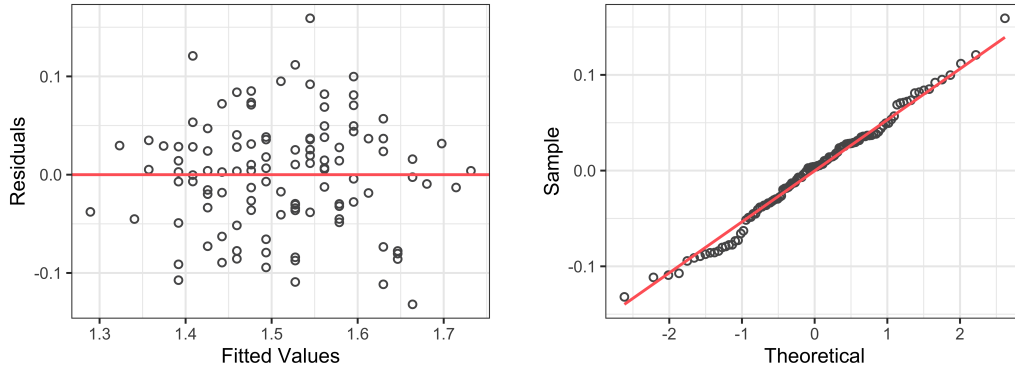


Figure 3: Diagnostic plots for the power transformation model.

- (e) We cannot conduct an  $F$ -test on the two models. This is because both models have the same number of parameters that need to be estimated, meaning they both have the same number of degrees of freedom.
- (f) In order to find a prediction interval for  $Y$ , we can first find a prediction interval for  $f_{0.085}(Y)$  and then apply  $f_{0.085}^{-1}$  to both ends of the interval. The prediction intervals for when  $X = 30$  and  $X = 60$  have been printed in Table 1. For reference, the prediction intervals at these observations for the linear model  $Y = \beta_0 + \beta_1 + \epsilon$  have also been printed. The prediction intervals for the power model can also be seen in the middle and right plots of Figure 2. For starters, the width of the intervals for the power model drastically increases as  $X$  increases, while the increase for the linear model is negligible. When  $X = 30$ , the confidence interval for the power model is significantly narrower than the linear function, but when  $X = 60$ , the opposite is true. That is, when  $X = 30$  we are confident that the next observation will lie in a smaller interval using the power model, but when  $X = 60$ , the linear model gives us the smaller region. Depending on the value of  $X$ , I would use the better of the two models.

	$X = 30$				$X = 60$			
Model	$\hat{Y}$	Lower	Upper	$\Delta\mathcal{I}$	$\hat{Y}$	Lower	Upper	$\Delta\mathcal{I}$
Linear	74.612	-89.781	239.004	328.785	297.769	133.036	462.503	329.467
Power	69.448	26.573	168.796	142.223	276.313	117.843	611.638	493.795

Table 1: Comparing the prediction intervals of the linear model against the power model.

## Question 2

*Collaborators:* None

## Appendix A

In this section, we will derive the result for  $m(\lambda)$ , the log-likelihood as a function of only  $\lambda$ . Suppose we have our response vector  $\mathbf{y}$  and our observed data  $\mathbf{X} = [\mathbf{1} \ \mathbf{x}]$ . We are interested in fitting the model  $\mathbf{g}_\lambda(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{g}_\lambda(\mathbf{y})$  is the transformed response vector (i.e. the  $i$ th element is given by  $g_\lambda(y_i)$ ),  $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ , and  $\text{Var}[\boldsymbol{\epsilon}] = \sigma^2\mathbf{I}$ . For notational ease, we will denote  $\mathbf{g}_\lambda(\mathbf{y})$  as  $\mathbf{g}_\lambda$ . If we make the further assumption that  $\boldsymbol{\epsilon}$  is normally distributed, i.e.  $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$ , then our response vector is also normally distributed, where  $\mathbf{g}_\lambda \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$ . It's density function (and thus it's likelihood function) is given by

$$\begin{aligned} f(\mathbf{g}_\lambda | \boldsymbol{\beta}, \sigma^2, \lambda) &= \frac{1}{\sqrt{\det(2\pi\sigma^2\mathbf{I})}} \cdot \exp\left(-\frac{(\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})^T (\sigma^2\mathbf{I})^{-1} (\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})}{2}\right) \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}\right). \end{aligned}$$

Since  $\mathbf{y}$  is a transformation of  $\mathbf{g}_\lambda$  (via  $g_\lambda^{-1}$ ), we can derive the density for  $\mathbf{y}$  as well. Notationally, this result may be somewhat confusing; even though we are finding the density for  $\mathbf{y}$ , we will still express the density (partly) in terms of  $\mathbf{g}_\lambda$ . It is important to remember that  $\mathbf{g}_\lambda$  is a function of  $\mathbf{y}$ . Because the  $i$ th element of  $\mathbf{g}_\lambda$  only depends on the  $i$ th element of  $\mathbf{y}$ , the Jacobian will be a diagonal matrix, and so

$$\mathbf{J} = \frac{\partial \mathbf{g}_\lambda}{\partial \mathbf{y}} = \text{diag}\left(\frac{\partial g_\lambda(y_1)}{\partial y_1}, \dots, \frac{\partial g_\lambda(y_n)}{\partial y_n}\right) = \text{diag}(y_1^{\lambda-1}, \dots, y_n^{\lambda-1}),$$

and so the density (and thus the likelihood) of  $\mathbf{y}$  is given by

$$g(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \lambda) = f(\mathbf{g}_\lambda(\mathbf{y})) \cdot |\det(\mathbf{J})| = \frac{1}{(2\pi\sigma^2)^{n/2}} \cdot \exp\left(-\frac{(\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}\right) \cdot \prod_{i=1}^n y_i^{\lambda-1}.$$

The log-likelihood  $\ell(\mathbf{y}) = \log g(\mathbf{y})$  is given by

$$\ell(\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \lambda) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{(\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{g}_\lambda - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2} + (\lambda - 1) \sum_{i=1}^n \log(y_i).$$

As is standard with maximum likelihood estimation, we now differentiate  $\ell$  with respect to the unknown parameters, set the derivatives to zero, and solve to get the maximum value of  $\ell$ . For now, we are going to leave  $\lambda$  fixed and differentiate with respect to  $\boldsymbol{\beta}$  and  $\sigma^2$ . Doing this for both gives us  $\hat{\boldsymbol{\beta}}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{g}_\lambda$  and  $\hat{\sigma}_{\text{MLE}}^2 = \mathbf{g}_\lambda^T (\mathbf{I} - \mathbf{H}) \mathbf{g}_\lambda / n$ , where  $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is the hat matrix. It is worth noting that both  $\hat{\boldsymbol{\beta}}_{\text{MLE}}$  and  $\hat{\sigma}_{\text{MLE}}^2$  are functions of  $\lambda$ . Plugging these values back into  $\ell$  will maximize it with respect to  $\boldsymbol{\beta}$  and  $\sigma^2$ , which means we will only have to maximize it with respect to  $\lambda$ . With some simplification, our new loss function is

$$m(\lambda) := \ell(\mathbf{y} | \hat{\boldsymbol{\beta}}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2, \lambda) = -\frac{n}{2} \log\left(\frac{2\pi e}{n}\right) - \frac{n}{2} \log(\mathbf{g}_\lambda^T (\mathbf{I} - \mathbf{H}) \mathbf{g}_\lambda) + (\lambda - 1) \sum_{i=1}^n \log(y_i).$$

Ideally, we would differentiate  $m$  with respect to  $\lambda$ , set  $\partial m / \partial \lambda = 0$ , and solve for  $\lambda$ . I was unable to derive a closed form solution for the result, but it is still possible to use graphical techniques or numerical methods to find the optimal value of  $\lambda$ .

We recall that both  $\hat{\beta}_{\text{MLE}}$  and  $\hat{\sigma}_{\text{MLE}}^2$  are functions of  $\lambda$ , so we cannot know their value until  $\hat{\lambda}_{\text{MLE}}$  has been determined. Because of this, as we just showed, the likelihood function can be expressed as a function  $m(\lambda)$  that only depends on  $\lambda$ , which can be maximized to find  $\hat{\lambda}_{\text{MLE}}$ . Once we find  $\hat{\lambda}_{\text{MLE}}$ , we can use this value to determine  $\hat{\beta}_{\text{MLE}}$  and  $\hat{\sigma}_{\text{MLE}}^2$ .