# Module 5: Peer Reviewed Assignment
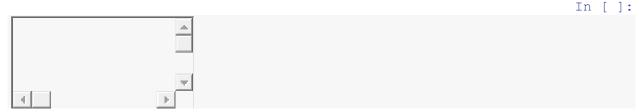
## Outline:

The objectives for this assignment:

1. Understand what can cause violations in the linear regression assumptions.
2. Enhance your skills in identifying and diagnosing violated assumptions.
3. Learn some basic methods of addressing violated assumptions.

General tips:

1. Read the questions carefully to understand what is being asked.
2. This work will be reviewed by another human, so make sure that you are clear and concise in what your explanations and answers.

In [ ]:

```
# Load Required Packages
library(ggplot2)
```
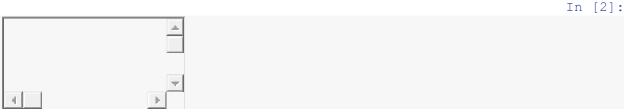
# Problem 1: Let's Violate Some Assumptions!

When looking at a single plot, it can be difficult to discern the different assumptions being violated. In the following problem, you will simulate data that purposefully violates each of the four linear regression assumptions. Then we can observe the different diagnostic plots for each of those assumptions.

**1. (a) Linearity**

Generate SLR data that violates the linearity assumption, but maintains the other assumptions. Create a scatterplot for these data using ggplot.

Then fit a linear model to these data and comment on where you can diagnose nonlinearity in the diagnostic plots.
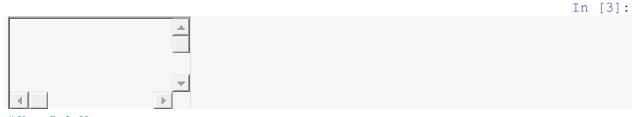
In [2]:

```
# Your Code Here
```

Type *Markdown* and LaTeX: $\alpha_2 \diamond 2$

**1. (b) Homoskedasticity**

Simulate another SLR dataset that violates the constant variance assumption, but maintains the other assumptions. Then fit a linear model to these data and comment on where you can diagnose non-constant variance in the diagnostic plots.
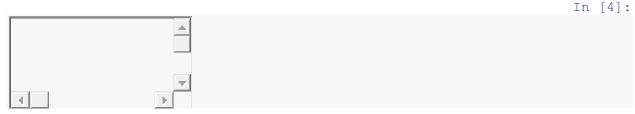
*# Your Code Here*

Type *Markdown* and LaTeX: $\alpha_2 \square^2$

**1. (c) Independent Errors**

Repeat the above process with simulated data that violates the independent errors assumption.

*# Your Code Here*

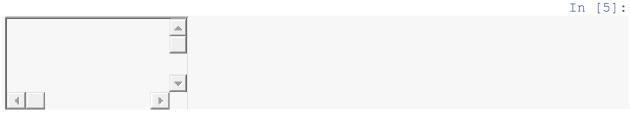Type *Markdown* and LaTeX: $\alpha_2 \square^2$

**1. (d) Normally Distributed Errors**

Only one more to go! Repeat the process again but simulate the data with non-normal errors.

*# Your Code Here*

Type *Markdown* and LaTeX: $\alpha_2 \square^2$

# Problem 2: Hats for Sale

Recall that the *hat* or *projection* matrix is defined as
$$H=(X_T X)_{-1}X_T. \square=\square(\square\square\square)^{-1}\square\square.$$
The goal of this question is to use the hat matrix to prove that the fitted values, $\mathbf{\hat{Y}}\square^\wedge$, and the residuals, $\hat{\varepsilon}\square^\wedge$, are uncorrelated. It's a bit of a process, so we will do it in steps.

**2. (a) Show that $\hat{Y}=HY\square^\wedge=\square\square$. That is, $H\square$ "puts a hat on" $Y\square$.**

Type *Markdown* and LaTeX: $\alpha_2 \square^2$

**2. (b) Show that $H\square$ is symmetric: $H=H_T\square=\square\square$.**

Type *Markdown* and LaTeX: $\alpha_2 \square 2$

**2. (c) Show that** $(I_n - H) = 0_n \square (\square\square - \square) = 0\square$**, where** $0_n 0\square$ **is the zero matrix of size** $n \times n \square \times \square$**.****

Type *Markdown* and LaTeX: $\alpha_2 \square 2$

**2. (d) Stating that** $\hat{Y}\square^{\wedge}$ **is uncorrelated with** $\hat{\varepsilon}\square^{\wedge}$ **is equivalent to showing that these vectors are orthogonal.* That is, we want their dot product to equal zero:**

$$\hat{Y}^T \hat{\varepsilon} = 0. \square^{\wedge}\square\square^{\wedge} = 0.$$

Prove this result. Also explain why being uncorrelated, in this case, is equivalent to the being orthogonal.

Type *Markdown* and LaTeX: $\alpha_2 \square 2$

**2.(e) Why is this result important in the practical use of linear regression?**

Type *Markdown* and LaTeX: $\alpha_2 \square 2$

# Problem 3: Model Diagnosis

We here at the University of Colorado's Department of Applied Math love Bollywood movies. So, let's analyze some data related to them!
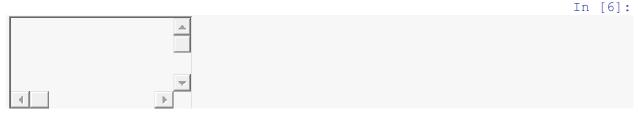
We want to determine if there is a linear relation between the amount of money spent on a movie (it's budget) and the amount of money the movie makes. Any venture capitalists among you will certianly hope that there is at least some relation. So let's get to modelling!

**3. (a) Initial Inspection**

Load in the data from local directory and create a linear model with `Gross` as the response and `Budget` as the feature. The data is stored in the same local directory and is called `bollywood_boxoffice.csv`. Thank the University of Florida for this specific dataset.

Specify whether each of the four regression model assumptions are being violated.

Data Source: http://www.bollymoviereviewz.com

In [6]:

```
# Load the data
bollywood = read.csv("bollywood_boxoffice.csv")
summary(bollywood)

# Your Code Here
```

```
        Movie              Gross              Budget
 1920London       : 1   Min.   :  0.63   Min.   :  4.00
 2 States\xa0      : 1   1st Qu.:  9.25   1st Qu.: 19.00
 24(Tamil,Telugu) : 1   Median : 29.38   Median : 34.50
 Aashiqui 2       : 1   Mean   : 53.39   Mean   : 45.25
```

```
AeDilHainMushkil\xa0:   1     3rd Qu.: 70.42     3rd Qu.: 70.00
AGentleman          :   1     Max.    :500.75    Max.    :150.00
(Other)             :184
```
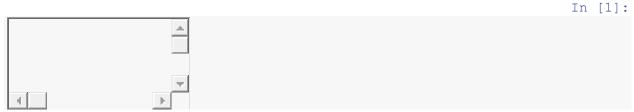
Type *Markdown* and LaTeX: $\alpha2\diamond2$

### 3. (b) Transformations

Notice that the Residuals vs. Fitted Values plot has a 'trumpet" shape to it, the points have a greater spread as the Fitted value increases. This means that there is not a constant variance, which violates the homoskedasticity assumption.

So how do we address this? Sometimes transforming the predictors or response can help stabilize the variance. Experiment with transfomrations on `Budget` and/or `Gross` so that, in the transformed scale, the relationship is approximately linear with a constant variance. Limit your transformations to square root, logarithms and exponentiation.

Note: There may be multiple transformations that fix this violation and give similar results. For the purposes of this problem, the transformed model doesn't have the be the "best" model, so long as it maintains both the linearity and homoskedasticity assumptions.

*# Your Code Here*

### 3. (c) Interpreting Your Transformation

You've fixed the nonconstant variance problem! Hurray! But now we have a transformed model, and it will have a different interpretation than a normal linear regression model. Write out the equation for your transformed model. Does this model have an interpretation similar to a standard linear model?

Type *Markdown* and LaTeX: $\alpha2$