# HW2 STA521

*[Your Name Here, netid and github username here]*

*Due September 12, 2019 10am*

## Background Reading

Readings: Chapters 3-4, 8-9 and Appendix in Weisberg Applied Linear Regression

This exercise involves the UN data set from `alr3` package. Install `alr3` and the `car` packages and load the data to answer the following questions adding your code in the code chunks. Please add appropriate code to the chunks to suppress messages and warnings as needed once you are sure the code is working properly and remove instructions if no longer needed. Figures should have informative captions. Please switch the output to pdf for your final version to upload to Sakai. **Remove these instructions for final submission**

## Exploratory Data Analysis

0. Preliminary read in the data. After testing, modify the code chunk so that output, messages and warnings are suppressed. *Exclude text from final*

```
library(alr3)
```

```
## Loading required package: car
```

```
## Loading required package: carData
```

```
data(UN3, package="alr3")
help(UN3)
library(car)
```

1. Create a summary of the data. How many variables have missing data? Which are quantitative and which are qualtitative?

2. Investigate the predictors graphically, using scatterplots or other tools of your choice. Create some plots highlighting the relationships among the predictors. Comment on your findings regarding trying to predict `ModernC` from the other variables. Are there potential outliers, nonlinear relationships or transformations that appear to be needed based on your graphical EDA?

## Model Fitting

3. Use the `lm()` function to perform a multiple linear regression with `ModernC` as the response and all other variables as the predictors, using the formula `ModernC ~ .`, where the `.` includes all remaining variables in the dataframe. Create diagnostic residual plot from the linear model object and comment on results regarding assumptions. How many observations are used in your model fitting?

4. Examine added variable plots `car::avPlot` or `car::avPlots` for your model above. Are there any plots that suggest that transformations are needed for any of the terms in the model? Describe. Is it likely that any of the localities are influential for any of the terms? Which localities? Which terms?

5. Using the multivariate BoxCox `car::powerTransform` find appropriate transformations of the response and predictor variables for the linear model. If any predictors are negative, you may need to transform so that they are non-negative. Summarize the resulting transformations.

6. Given the selected transformations of the predictors, verify the transformation of the response using `MASS::boxcox` or `car::boxCox` and justify. Do you get the same transformation if you used `car::powerTransform` above? Do you get the same transformation for the response if you do not transform any of the predictors? Discuss briefly the findings.

7. Fit the regression using the transformed variables. Provide residual plots and added variables plots and comment. If you feel that you need additional transformations of either the response or predictors, repeat any steps until you feel satisfied with the model and residuals.

8. Are there any outliers or influential points in the data? Explain. If so, refit the model after removing any outliers/influential points and comment on residual plots.

## Summary of Results

9. For your final model, provide summaries of coefficients with 95% confidence intervals in a nice table with interpretations of each coefficient. These should be in terms of the original units!

10. Provide a paragraph summarizing your final model and findings suitable for the US envoy to the UN after adjusting for outliers or influential points. You should provide a justification for any case deletions in your final model.

## Methodology

11. Exercise 9.12 from ALR

Using $X^T X = X_{(i)}^T X_{(i)} + x_i x_i^T$ where the subscript $(i)$ means without the ith case, show that

$$(X_{(i)}^T X_{(i)})^{-1} = (X^T X)^{-1} + \frac{(X^T X)^{-1} x_i x_i^T (X^T X)^{-1}}{1 - h_{ii}}$$

where $h_{ii}$ is the $i$th diagonal element of $H = X(X^T X)^{-1} X^T$ using direct multiplication and simplify in terms of $h_{ii}$.

12. Exercise 9.13 from ALR. Using the above, show

$$\hat{\beta}_{(i)} = \hat{\beta} - \frac{(X^T X)^{-1} x_i e_i}{1 - h_{ii}}$$

13. (optional) Prove that the intercept in the added variable scatter plot will always be zero. *Hint: use the fact that if $H$ is the projection matrix for $X$ which contains a column of ones, then $1_n^T(I - H) = 0$ or $(I - H)1_n = 0$. Use this to show that the sample mean of residuals will always be zero if there is an intercept.*