# 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

```r
#install.packages("alr3")
#install.packages("car")
#install.packages("backports")
#install.packages("readxl")
#install.packages("htmltools")
#install.packages('scales')
#install.packages("colorspace")
#install.packages("lazyeval")
#install.packages("yaml")
#install.packages("stringi")
```

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*

```r
library(alr3)
```

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

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

```r
data(UN3, package="alr3")
help(UN3)
```

```
## starting httpd help server ...
```

```
##  done
```

```r
library(car)
```

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

```r
str(UN3)
```

```
## 'data.frame':    210 obs. of  7 variables:
##  $ ModernC  : int  NA NA 49 NA NA NA 51 NA 22 NA ...
##  $ Change   : num  3.88 0.68 1.67 2.37 2.59 3.2 0.53 1.17 -0.45 2.02 ...
##  $ PPgdp    : int  98 1317 1784 NA 14234 739 8461 7163 687 NA ...
##  $ Frate    : int  NA NA 7 42 NA NA 63 44 51 53 ...
##  $ Pop      : num  23897 3167 31800 57 64 ...
##  $ Fertility: num  6.8 2.28 2.8 NA NA 7.2 NA 2.44 1.15 NA ...
##  $ Purban   : int  22 43 58 53 92 35 37 88 67 51 ...
```

```r
smry_un3 <- summary(UN3)
na_count <- smry_un3[7,]
na_count
```

```
##       ModernC          Change          PPgdp           Frate            Pop
## "NA's   :58 "  "NA's   :1  "  "NA's   :9  "  "NA's   :43 "  "NA's   :2  "
##     Fertility          Purban
## "NA's   :10 "             NA
```

answer: As we can see in outlook of data.frame UN3, there are all quantative variables. Except for variable named Purban, those of variables including ModernC,Change,PPgdp,Frate,Pop,Fertility have at least one missing data.

2. What is the mean and standard deviation of each quantitative predictor? Provide in a nicely formatted table.

```r
library(knitr)
mn_st_table <- matrix(rep(0,3*length(UN3)),nrow = length(UN3))
for(i in 1:length(UN3)){
  mn_st_table[i,] <- c(colnames(UN3)[i],
                       round(mean(UN3[,i],na.rm = T),3),
                       round(sd(UN3[,i],na.rm=T),3))
}
rownames(mn_st_table) <- 1:length(UN3)
colnames(mn_st_table) <- c("variable","mean","stand deviation")
kable(mn_st_table)
```
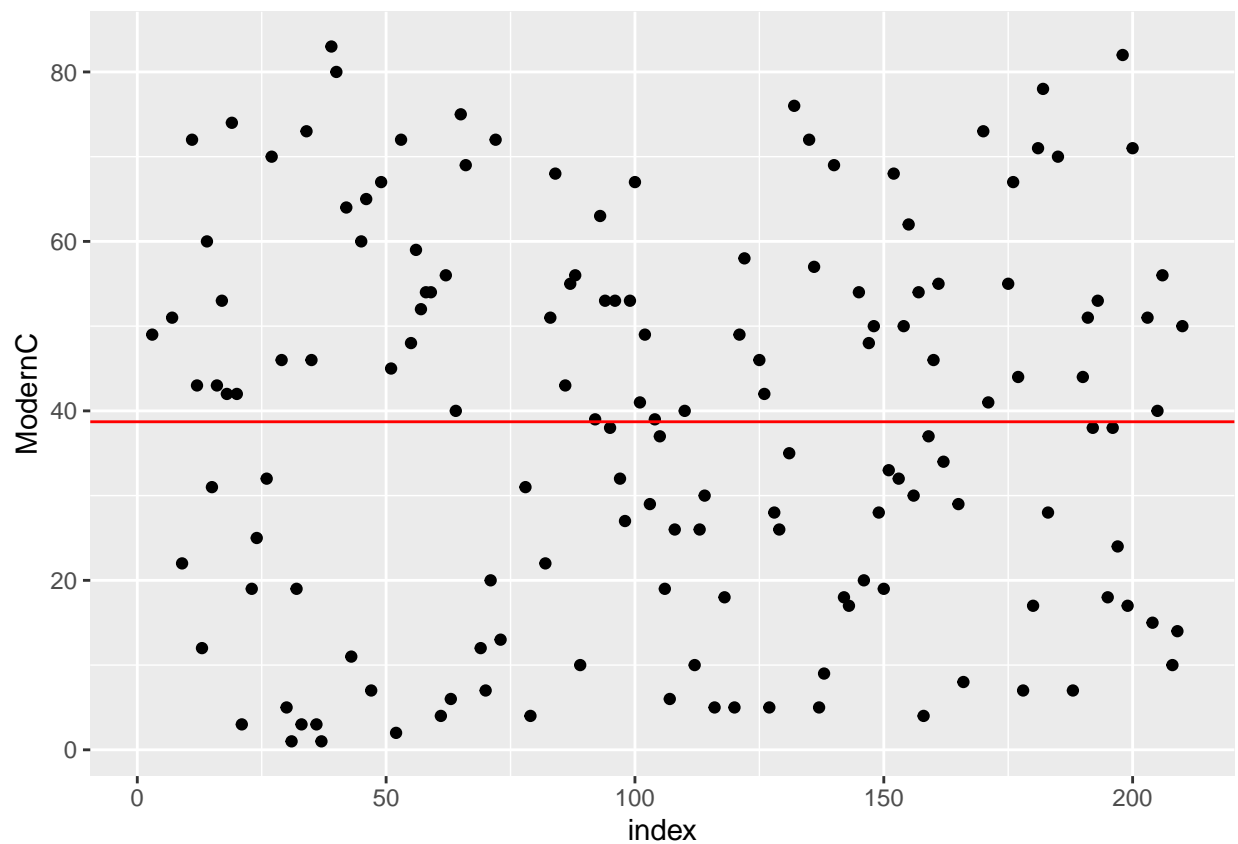
| variable | mean | stand deviation |
|----------|------|-----------------|
| ModernC | 38.717 | 22.637 |
| Change | 1.418 | 1.133 |
| PPgdp | 6527.388 | 9325.189 |
| Frate | 48.305 | 16.532 |
| Pop | 30281.871 | 120676.694 |
| Fertility | 3.214 | 1.707 |
| Purban | 56.2 | 24.11 |

3. 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
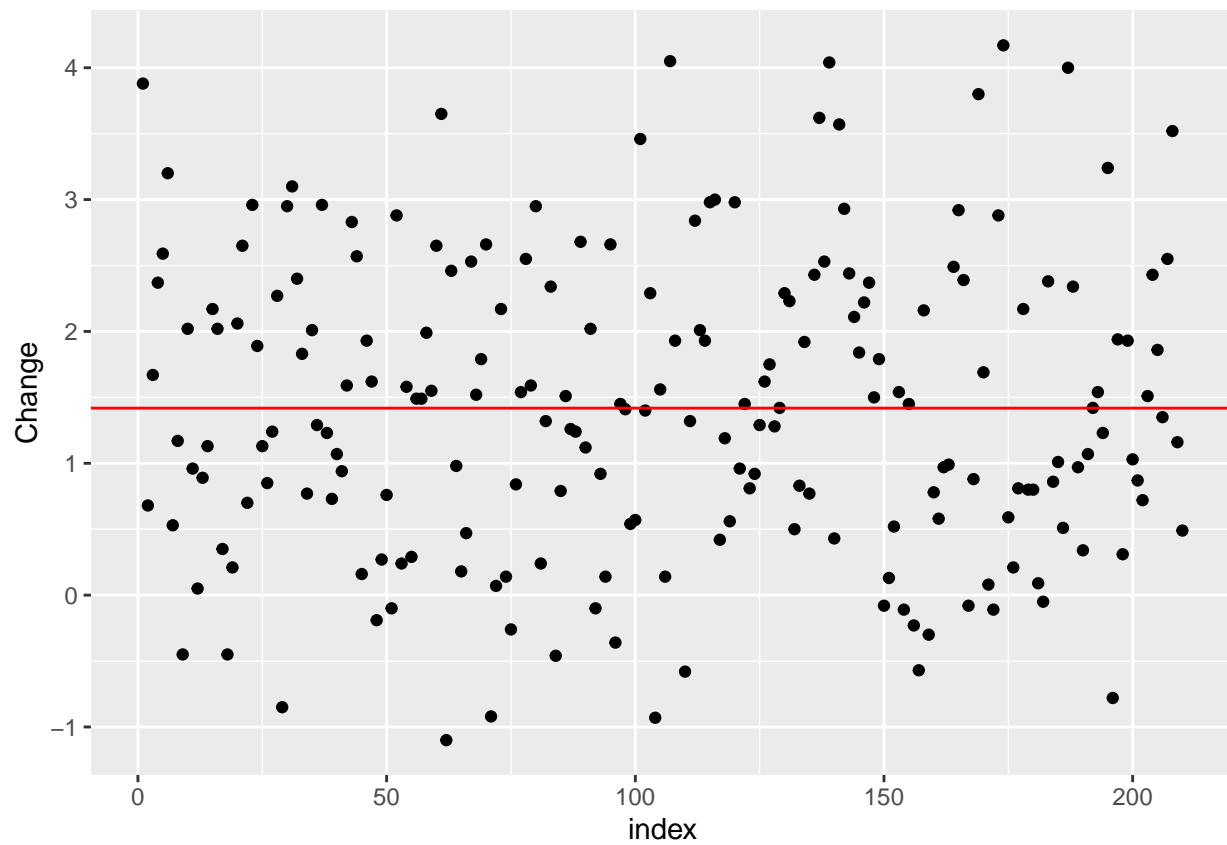
to predict `ModernC` from the other variables. Are there potential outliers, nonlinear relationships or transformations that appear to be needed based on your graphical EDA?

```
library(ggplot2)
par(mfrow=c(3,3))
for(i in 1:length(UN3)){
  print(ggplot(data=UN3, mapping = aes(x=1:nrow(UN3),y=UN3[,i])) +
    geom_point() +
    geom_hline(yintercept = mean(UN3[,i],na.rm = T),color="red") +
    ylab(colnames(UN3)[i]) + xlab("index"))
}
```
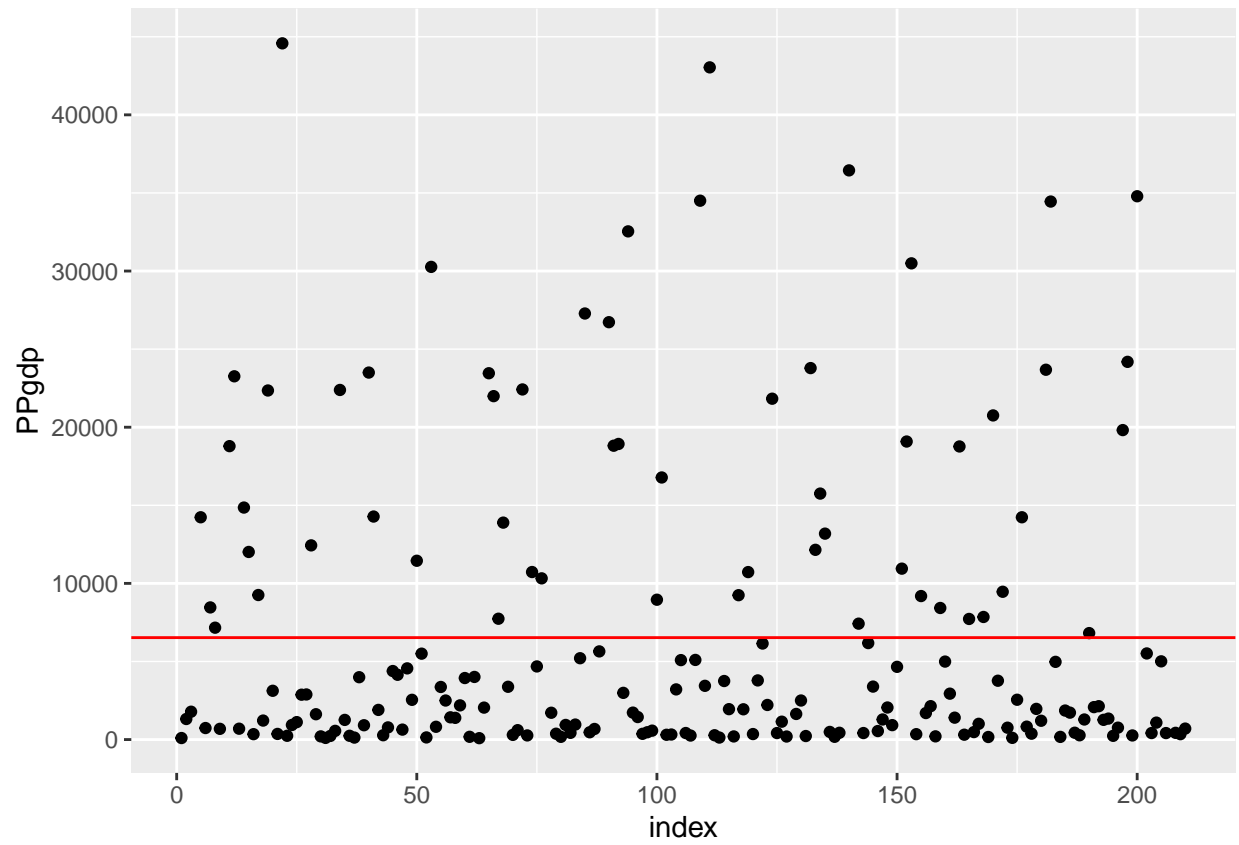
## Warning: Removed 58 rows containing missing values (geom_point).



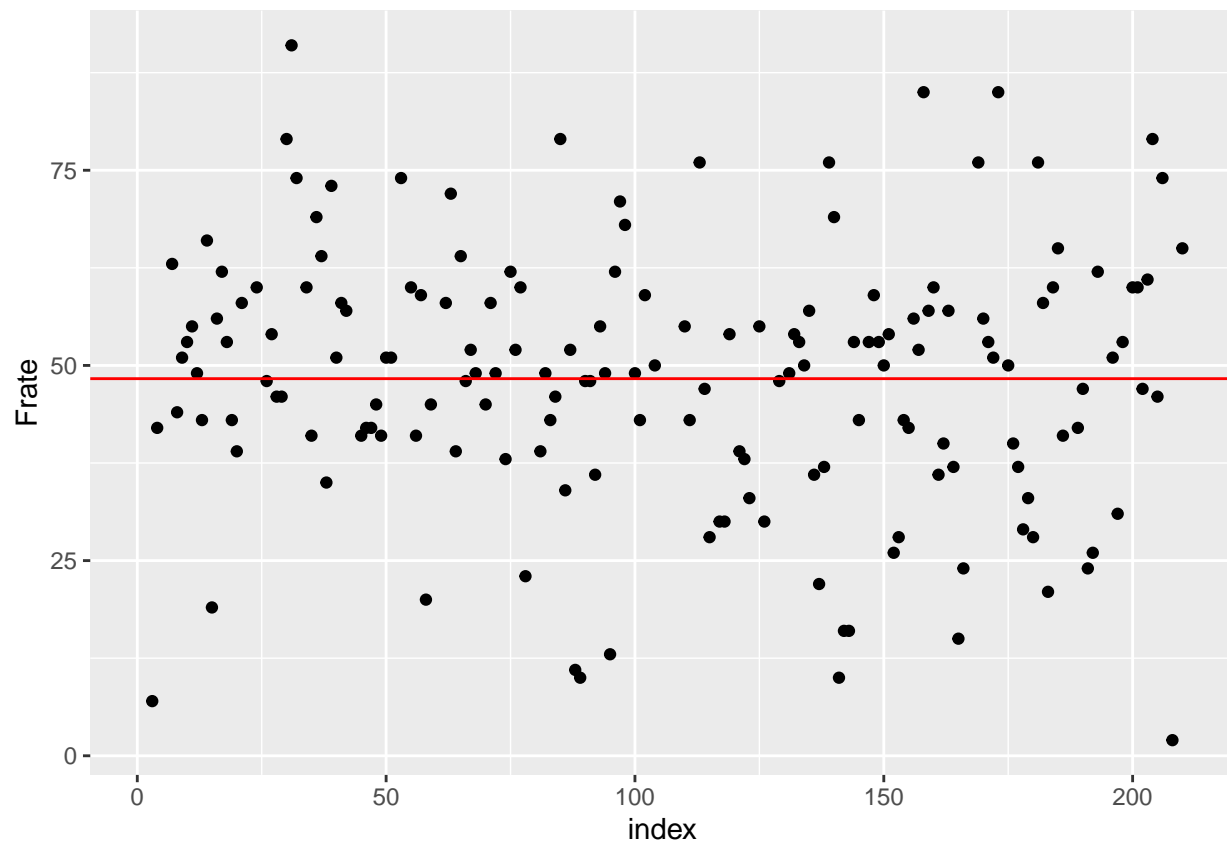## Warning: Removed 1 rows containing missing values (geom_point).

## Warning: Removed 9 rows containing missing values (geom_point).
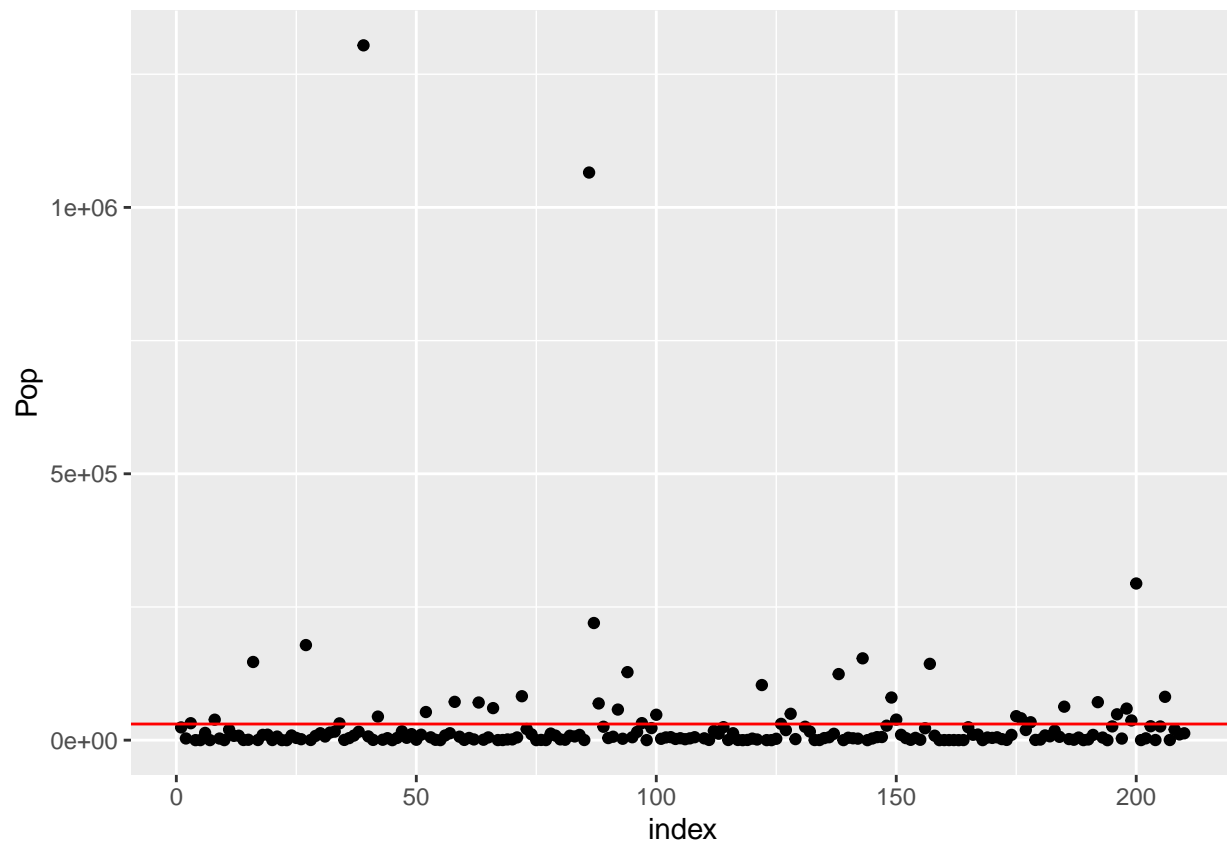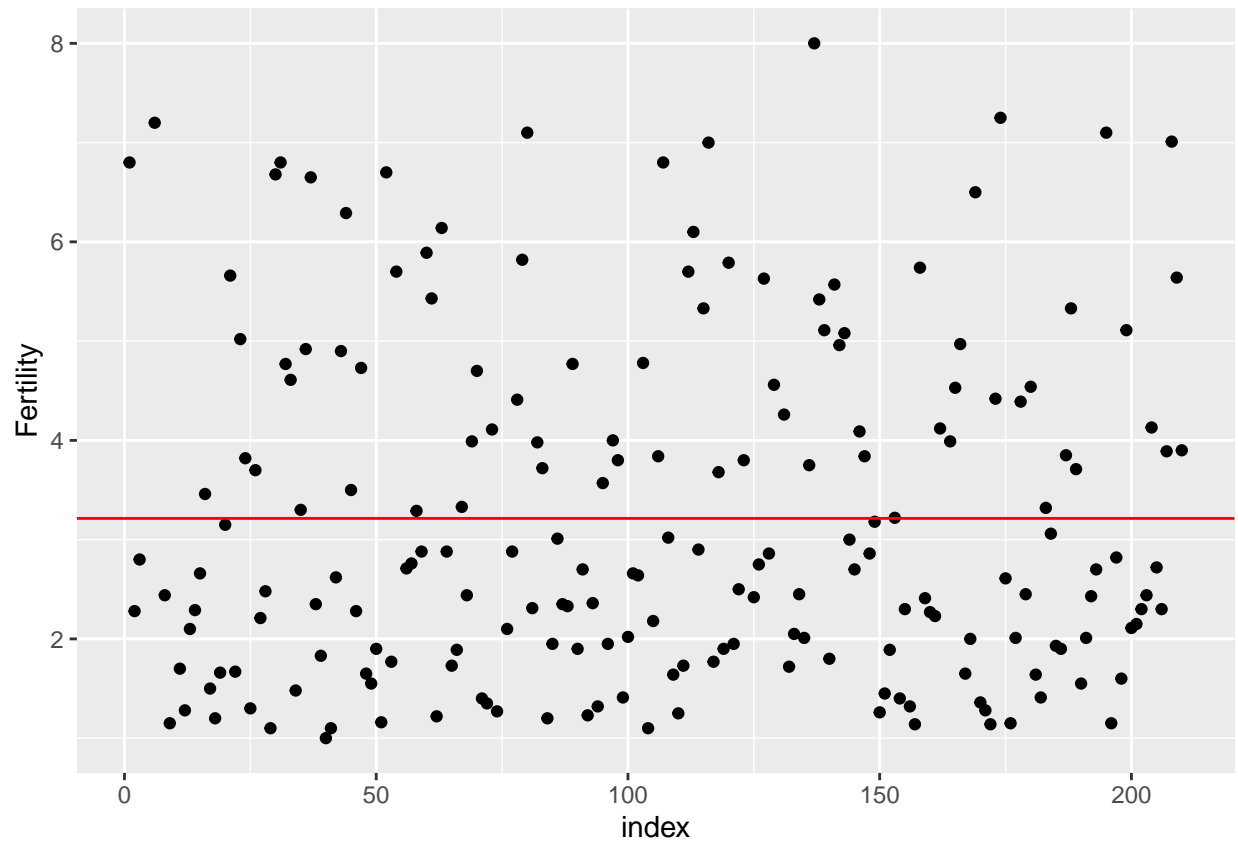
```
## Warning: Removed 43 rows containing missing values (geom_point).
```
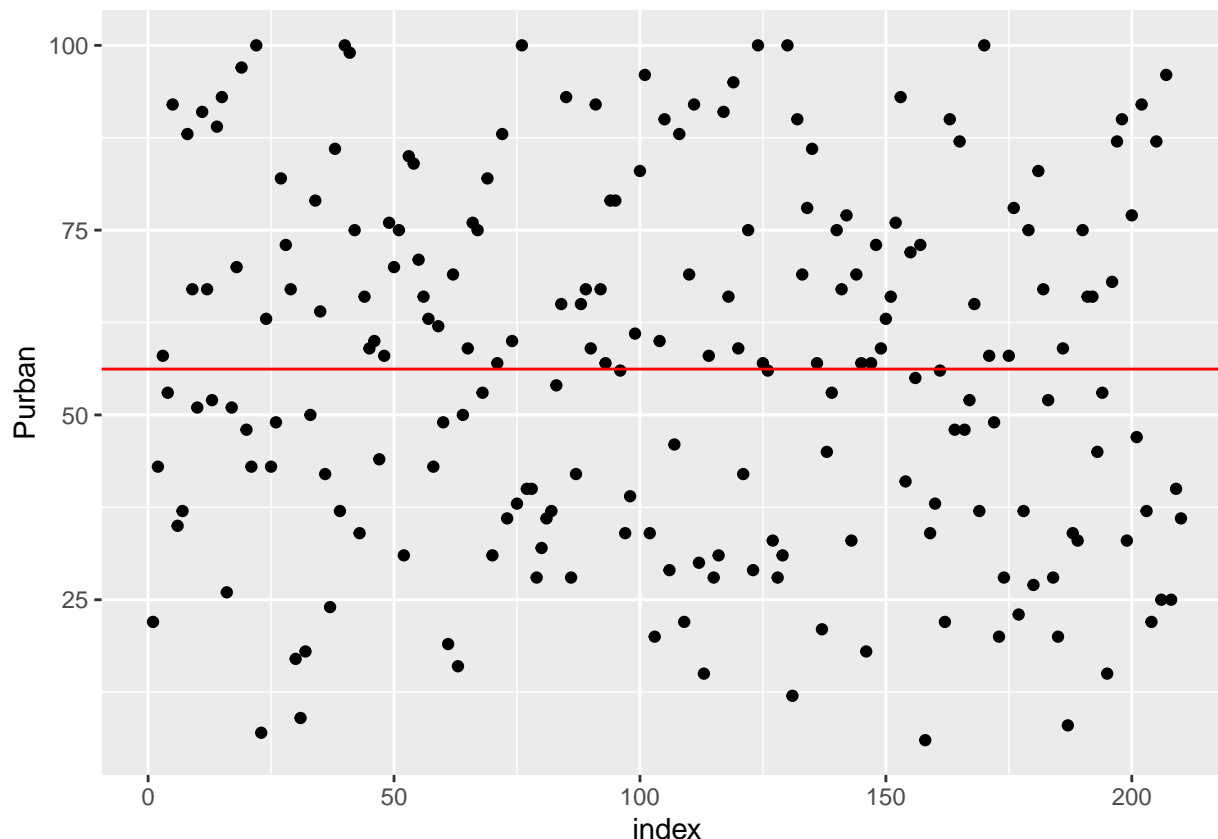
```
## Warning: Removed 2 rows containing missing values (geom_point).
```

```
## Warning: Removed 10 rows containing missing values (geom_point).
```

## Model Fitting

4. 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?

5. 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?

6. Using the multivariate BoxCox `car::powerTransform` or Box-Tidwell `car::boxTidwell` find appropriate transformations of the predictor variables to be used as predictors in the linear model. If any predictors are negative, you may need to transform so that they are non-negative. Describe your method and the resulting transformations.

7. Given the selected transformations of the predictors, select a transformation of the response using `MASS::boxcox` or `car::boxCox` and justify. Do you need to do this if you used `car::powerTransform` above? Explain.

8. 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.

9. 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

10. 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!

11. 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

12. 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.*

13. 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__ $h_{ii}$.

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

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