

# HW 02

Group 10

2022-10-03

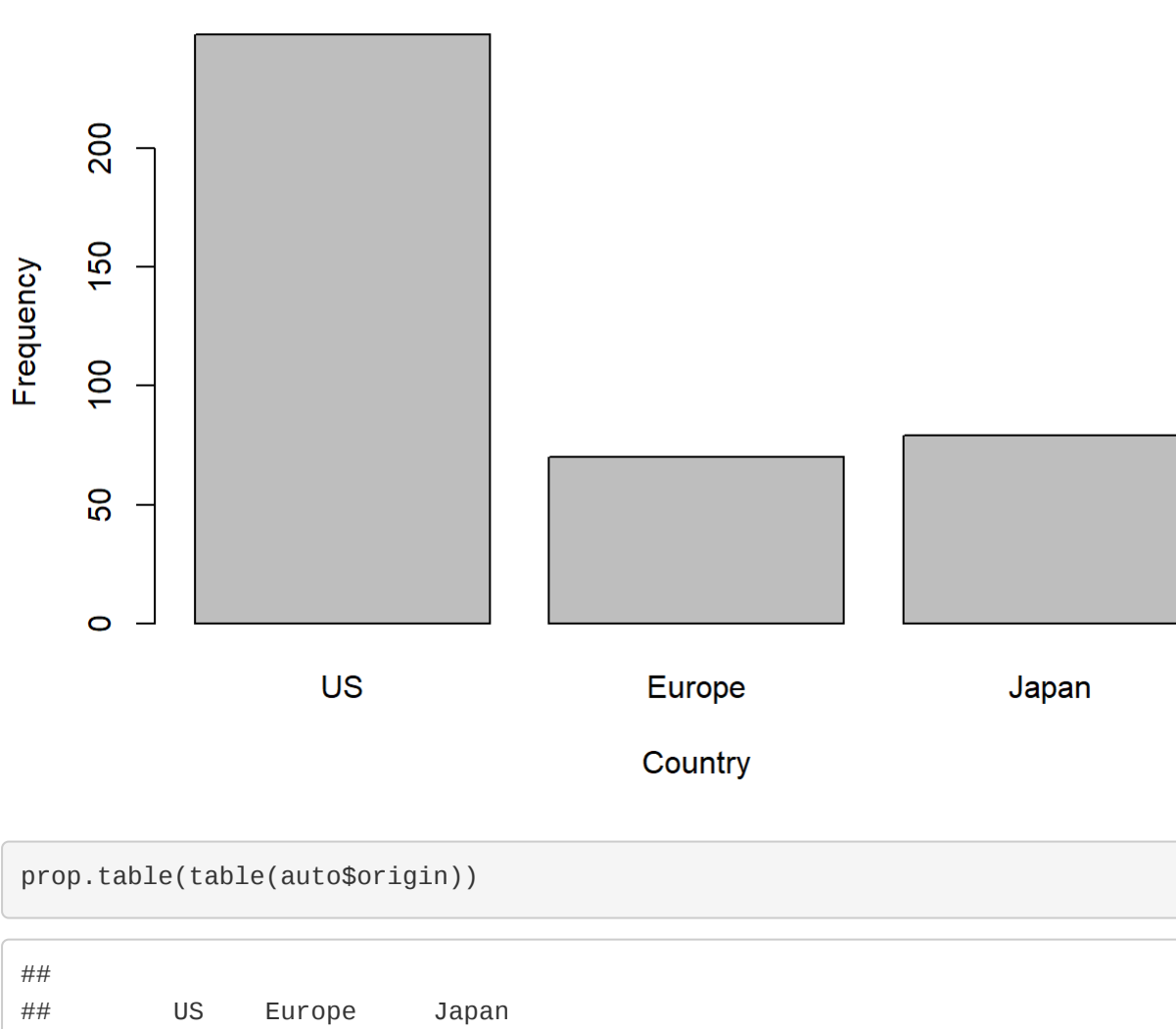
## Homework 2

Read csv

```
auto = read.csv("Auto.csv", na.strings = "NA")
```

1(a)

```
auto$origin = factor(auto$origin, 1:3, c("US", "Europe", "Japan"))  
freq <- table(auto$origin)  
barplot(freq, main = "Frequency of vehicle production in different countries", xlab = "Country", ylab = "Frequency")
```



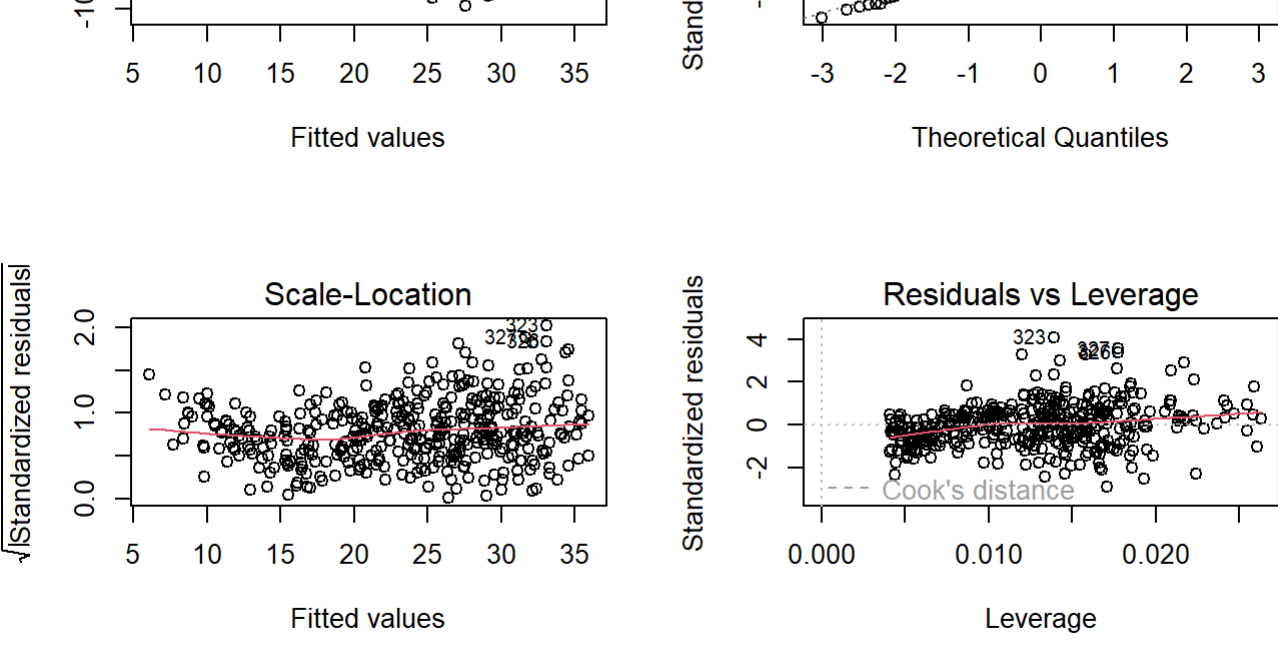
```
prop.table(table(auto$origin))
```

```
##  
##      US      Europe      Japan  
## 0.6246851 0.1763224 0.1989924
```

Above you can see the frequency is much greater in the US.

1(b)

```
lmb <- lm(mpg ~ origin + weight + year, data = auto)  
par(mfrow=c(2,2))  
plot(lmb)
```



The plots indicate that: - The error term is not normally distributed - The variance is not constant, thus heteroskedasticity is present - The data is not normally distributed

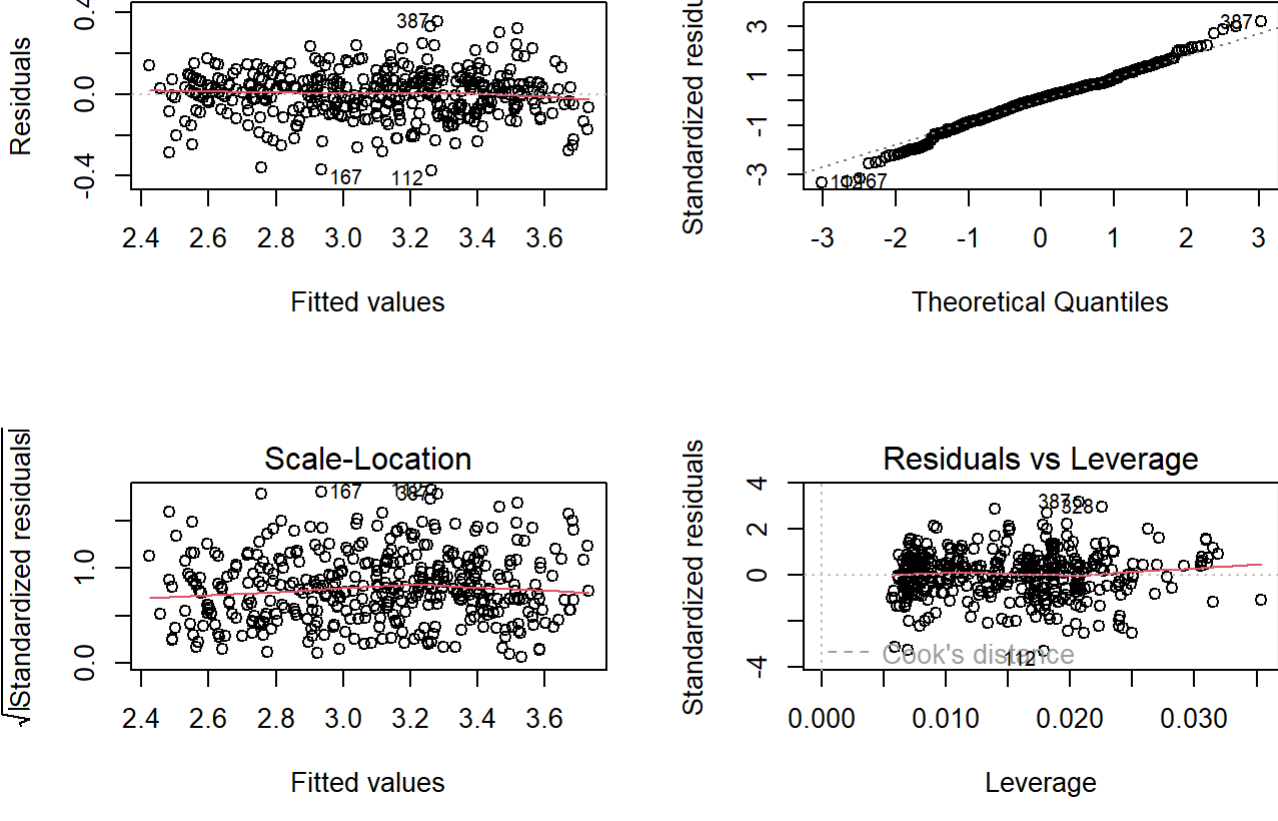
We can tell that the error term is not normally distributed by looking at the Residuals vs Fitted plot. The red line doesn't move along 0 at all. It looks more like a quadratic function.

The variance as x increases on the Residuals vs Fitted plot which is an indicator that heteroskedasticity is present.

We can see from the QQ plot that the data doesn't follow the line. Towards the top the data skews upwards.

1(c)

```
lmc <- lm(log(mpg) ~ origin + log(weight) + year + I(year^2), data = auto)  
par(mfrow=c(2,2))  
plot(lmc)
```



```
summary(lmc)
```

```
##  
## Call:  
## lm(formula = log(mpg) ~ origin + log(weight) + year + I(year^2),  
## data = auto)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.37408 -0.06782  0.00899  0.06903  0.35766   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  18.4693014   2.6833895   6.883 2.34e-11 ***  
## originEurope  0.0668291   0.0176293   3.791 0.000174 ***  
## originJapan   0.0319711   0.0179382   1.782 0.075477 .  
## log(weight)   -0.8750305   0.0270390 -32.362 < 2e-16 ***  
## year          -0.2559684   0.0712094  -3.595 0.000366 ***  
## I(year^2)      0.0019051   0.0004687   4.065 5.81e-05 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1136 on 391 degrees of freedom  
## Multiple R-squared:  0.8898, Adjusted R-squared:  0.8884   
## F-statistic: 631.7 on 5 and 391 DF,  p-value: < 2.2e-16
```

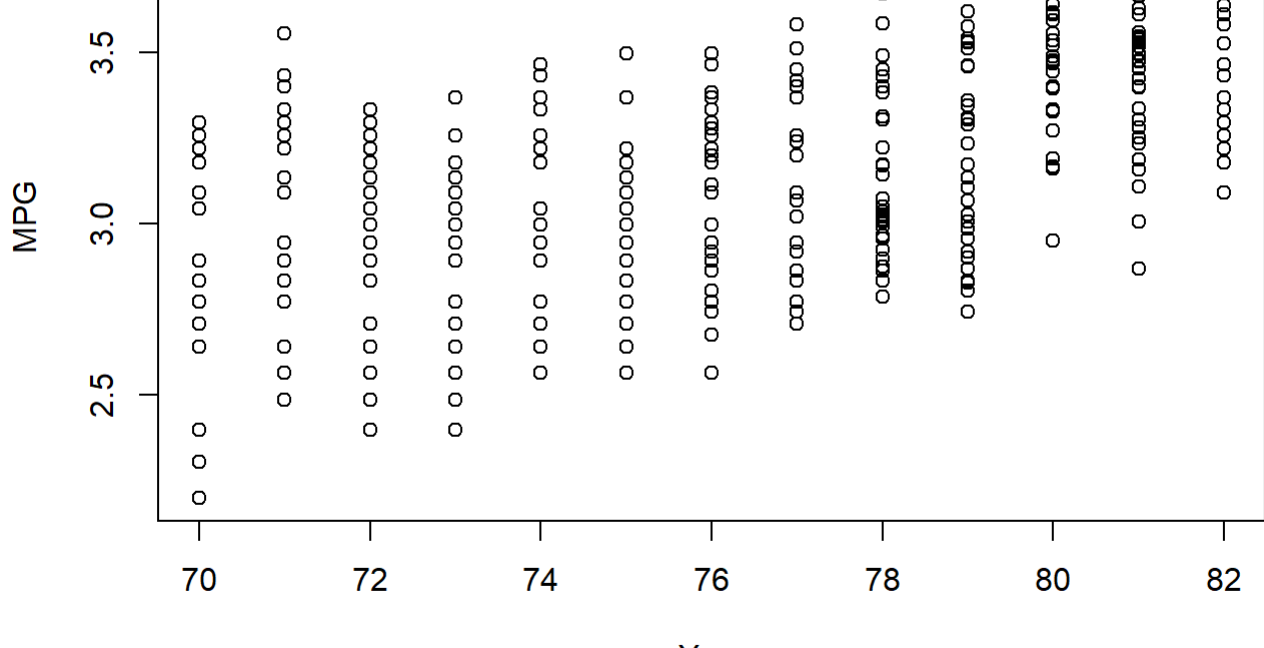
The model assumptions seem to have been roughly satisfied now.

The previously unsatisfied assumptions: - Heteroskedasticity - Error term is not normally distributed - Data is not normally distributed

When looking at the Residuals vs Fitted plot, we see the line follows 0 well and the variance is pretty much constant for each x value. We also see from the QQ plot that the data is more normally distributed now.

1(d)

```
plot(auto$year, log(auto$mpg),  
      main="Log(mpg) vs year",  
      xlab = "Year", ylab="MPG")
```



```
min = -coef(lmc)[5] / (2 * coef(lmc)[6])
```

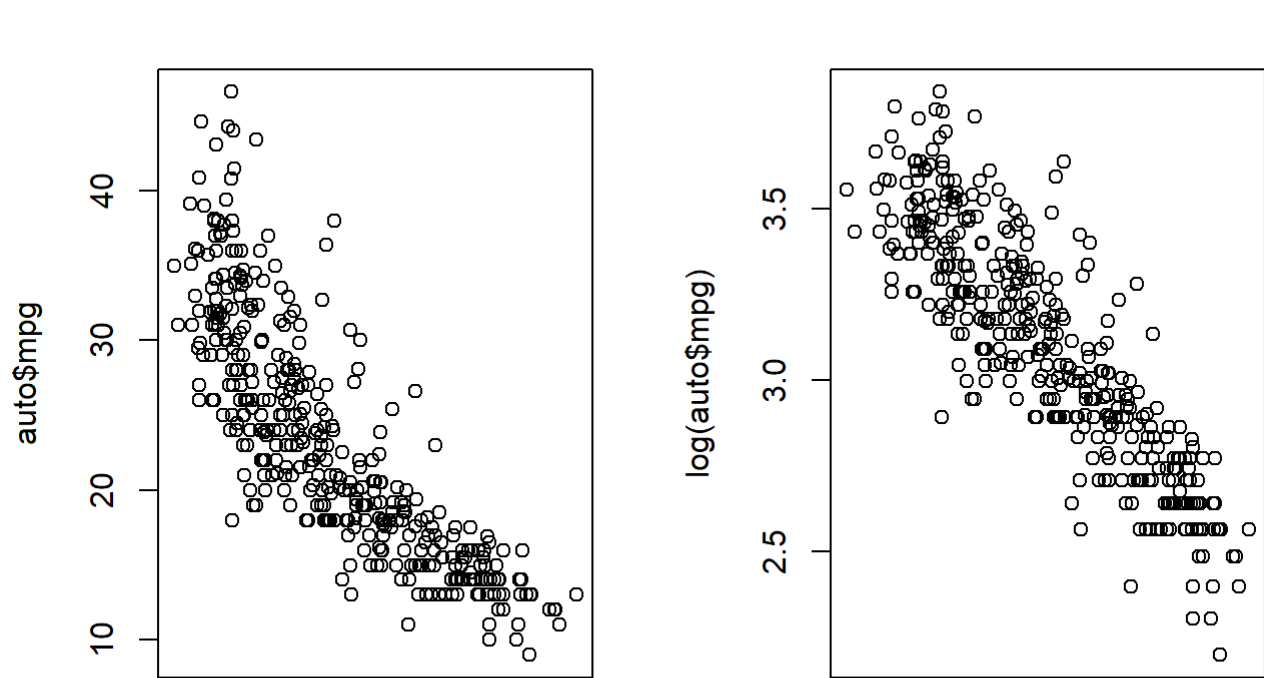
The relationship appears U-shaped based on the plot above. The minimum is 67.1781178.

1(e)

```
summary(lmc)
```

```
##  
## Call:  
## lm(formula = log(mpg) ~ origin + log(weight) + year + I(year^2),  
## data = auto)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.37408 -0.06782  0.00899  0.06903  0.35766   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  18.4693014   2.6833895   6.883 2.34e-11 ***  
## originEurope  0.0668291   0.0176293   3.791 0.000174 ***  
## originJapan   0.0319711   0.0179382   1.782 0.075477 .  
## log(weight)   -0.8750305   0.0270390 -32.362 < 2e-16 ***  
## year          -0.2559684   0.0712094  -3.595 0.000366 ***  
## I(year^2)      0.0019051   0.0004687   4.065 5.81e-05 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1136 on 391 degrees of freedom  
## Multiple R-squared:  0.8898, Adjusted R-squared:  0.8884   
## F-statistic: 631.7 on 5 and 391 DF,  p-value: < 2.2e-16
```

```
par(mfrow = c(1, 2))  
plot(auto$weight, auto$mpg)  
plot(log(auto$weight), log(auto$mpg))
```



It tells us that as you increase the weight the mpg falls. The relationship for the unlogged version is similar, less linear, but still negative.

2(a)

$$\begin{aligned}y_i &= \gamma_0 + \gamma_1(x_i - \bar{x}) + \gamma_2(x_i - \bar{x})^2 + e_i = \\&= \gamma_0 + \gamma_1 x_i - \gamma_1 \bar{x} + \gamma_2 x_i^2 - 2\gamma_2 x_i \bar{x} + \gamma_2 \bar{x}^2 + e_i = \\&= (\gamma_0 - \gamma_1 \bar{x} + \gamma_2 \bar{x}^2) + (\gamma_1 - 2\gamma_2 \bar{x})x_i + \gamma_2 x_i^2 + e_i \\&\therefore \beta_0 = \gamma_0 - \gamma_1 \bar{x} + \gamma_2 \bar{x}^2 \\&\quad \beta_1 = \gamma_1 - 2\gamma_2 \bar{x} \\&\quad \beta_2 = \gamma_2\end{aligned}$$

2(b)

```
auto$year_squared <- auto$year^2  
summary(lm(mpg ~ year + year_squared, data = auto))
```

```
##  
## Call:  
## lm(formula = mpg ~ year + year_squared, data = auto)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -13.349  -5.109  -0.878   4.587  18.196   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  577.25230   146.67144   3.936 9.81e-05 ***  
## year         -15.84090   3.86508  -4.098 5.05e-05 ***  
## year_squared  0.11230    0.02542  4.419 1.29e-05 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 6.23 on 394 degrees of freedom  
## Multiple R-squared:  0.3694, Adjusted R-squared:  0.3662   
## F-statistic: 115.4 on 2 and 394 DF,  p-value: < 2.2e-16
```

2(c)

The correlation between year and year squared is 0.999759.

2(d)

The mean of year is 75.9949622.

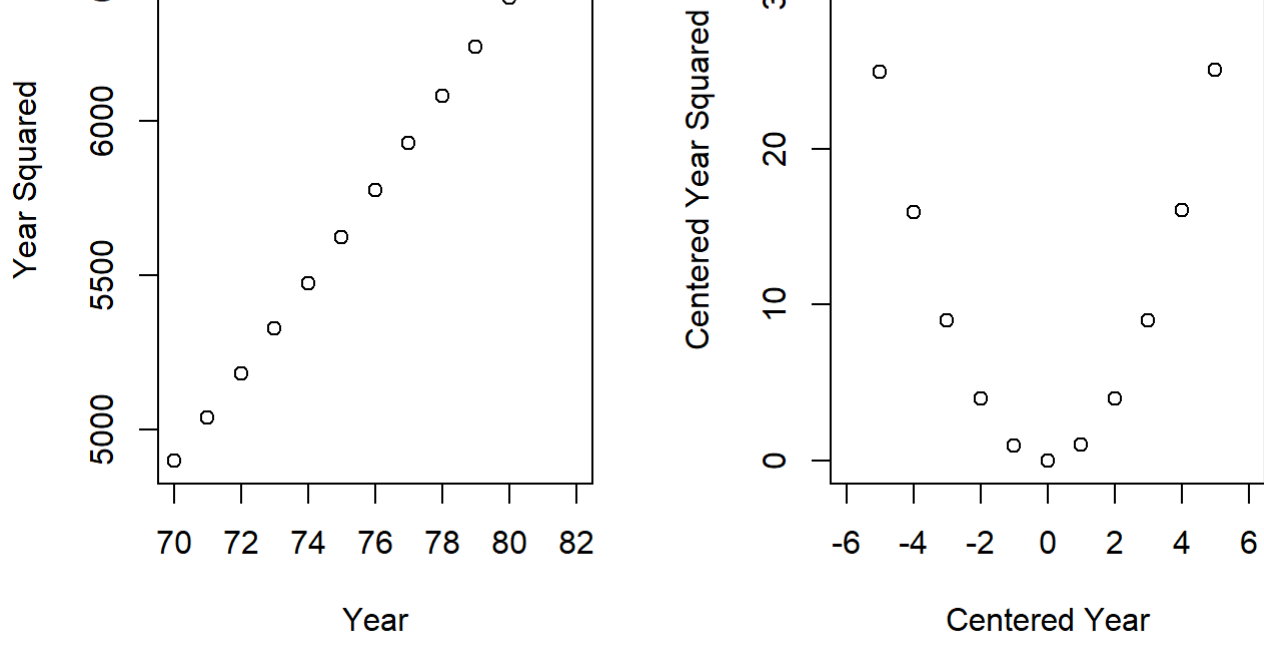
2(e)

```
auto$year_centered <- auto$year - mean(auto$year)  
auto$year_centered_squared <- (auto$year_centered)^2
```

The correlation between centered year and centered year squared is 0.014414

2(f)

```
par(mfrow = c(1, 2))  
plot(auto$year_centered, auto$year_centered_squared, main = "(Year vs Year Squared)", xlab = "Year", ylab = "Year Squared")  
plot(auto$year_centered, auto$year_centered_squared, main = "Centered (Year vs Year Squared)", xlab = "Centered Year", ylab = "Centered Year Squared")
```



2(g)

```
lm_2g <- lm(mpg ~ year_centered + year_centered_squared, data = auto)  
summary(lm_2g)
```

```
##  
## Call:  
## lm(formula = mpg ~ year_centered + year_centered_squared, data = auto)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -13.349  -5.109  -0.878   4.587  18.196   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  21.99061    0.46577  47.214 < 2e-16 ***  
## year_centered  1.22778    0.08486  14.469 < 2e-16 ***  
## year_centered_squared 0.11230    0.02542  4.419 1.29e-05 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 6.23 on 394 degrees of freedom  
## Multiple R-squared:  0.3694, Adjusted R-squared:  0.3662   
## F-statistic: 115.4 on 2 and 394 DF,  p-value: < 2.2e-16
```

2(h)

$$\begin{aligned}\beta_0 &= \gamma_0 - \gamma_1 \bar{x} + \gamma_2 \bar{x}^2 \\ \beta_1 &= \gamma_1 - 2\gamma_2 \bar{x} \\ \beta_2 &= \gamma_2\end{aligned}$$

```
gamma_0 <- coef(lm_2g)[1]  
gamma_1 <- coef(lm_2g)[2]  
gamma_2 <- coef(lm_2g)[3]  
mean_year <- mean(auto$year)
```

$$\begin{aligned}\beta_0 &= 577.2522975 \\ \beta_1 &= -15.8409008 \\ \beta_2 &= 0.1123014\end{aligned}$$