

Hw03ST430Yu

Haozhe (Jerry) Yu

2023-09-21

Question 1

```
educ <- as_tibble(read.table("https://users.stat.ufl.edu/~rrandles/sta4210/Rclassnotes/data/textdataset.
  #sep = "",
  strip.white=TRUE,
  col.name = c("Crime.Rate", "High.School.Diploma")
))
educ
```

```
## # A tibble: 84 x 2
##   Crime.Rate High.School.Diploma
##   <int>      <int>
## 1      8487         74
## 2      8179         82
## 3      8362         81
## 4      8220         81
## 5      6246         87
## 6      9100         66
## 7      6561         68
## 8      5873         81
## 9      7993         74
## 10     7932         82
## # i 74 more rows
```

a. Find the least squares regression equation to predict the crime rate from the percent of individuals having at least a high school education. [Paste R or SAS output and then answer your question]

```
educm <- lm(Crime.Rate~High.School.Diploma,data=educ)
```

The equation to predict crime rate (per 100,000 residents) from the percent of individuals in a country with at least a high school diploma is

Crime Rate = $2.05176 \times 10^4 + -170.5751886$ High School Percent

b. Give the ANOVA Table for this regression analysis. [Paste R or SAS output]

```
educma <- anova(educm)
```

```
educma
```

```
## Analysis of Variance Table
```

```
##
```

```
## Response: Crime.Rate
```

```
##           Df      Sum Sq Mean Sq F value    Pr(>F)
## High.School.Diploma  1  93462942 93462942  16.834 9.571e-05 ***
## Residuals          82  455273165  5552112
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
educma$"Pr(>F)"[1]
```

```
## [1] 9.571396e-05
```

```
educma$"F value"[1]
```

```
## [1] 16.83376
```

```
-sqrt(educma$"Sum Sq"[1]/(educma$"Sum Sq"[1] + educma$"Sum Sq"[2]))
```

```
## [1] -0.4127033
```

```
summary(educm)
```

```
##
```

```
## Call:
```

```
## lm(formula = Crime.Rate ~ High.School.Diploma, data = educ)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -5278.3 -1757.5  -210.5   1575.3   6803.3
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    20517.60    3277.64   6.260 1.67e-08 ***
## High.School.Diploma  -170.58     41.57  -4.103 9.57e-05 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 2356 on 82 degrees of freedom
```

```
## Multiple R-squared:  0.1703, Adjusted R-squared:  0.1602
```

```
## F-statistic: 16.83 on 1 and 82 DF, p-value: 9.571e-05
```

```
cor(educ$Crime.Rate,educ$High.School.Diploma)
```

```
## [1] -0.4127033
```

```
str(educma)
```

```
## Classes 'anova' and 'data.frame':  2 obs. of  5 variables:
## $ Df      : int  1 82
## $ Sum Sq : num  9.35e+07 4.55e+08
## $ Mean Sq: num  93462942 5552112
## $ F value: num  16.8 NA
## $ Pr(>F) : num  9.57e-05 NA
## - attr(*, "heading")= chr [1:2] "Analysis of Variance Table\n" "Response: Crime.Rate"
```

```
educma$"Sum Sq"[2]
```

```
## [1] 455273165
```

```
educma$"Sum Sq"[2]
```

```
## [1] 455273165
```

c. Find SSE and MSE for this model.

The SSE for this model is 4.5527317×10^8 and the MSE is 5.5521118×10^6

d. What is the estimate of sigma from this analysis?

The estimate of σ for this analysis is 2356.2919539

e. What percent of the variation in crime rates can be explained by the percent of high school graduates?

The percent of variation in crime rates explained by the percent of high school grads is 0.170324

f. What is the correlation between crime rates and percent of high school graduates?

The correlation between crime rates and percent of high school graduates is -0.4127033

g. Based on your ANOVA table, is the linear relationship between X and Y statistically significant? Be sure to give an appropriate null and alternate hypothesis, test statistic, its associated degrees of freedom, and the p-value.

- H_0 : There is no linear relationship between crime rates and percent of high school graduates ($\beta_1 = 0$)

- HA: There is a linear relationship between crime rates and percent of high school graduates ($\beta_1 \neq 0$)
- Test Statistic (F value): 16.8337645
- Degrees of Freedom: 1 for the model (High School Diploma Percent), and 82 for the error.
- P value: 9.5713958×10^{-5}

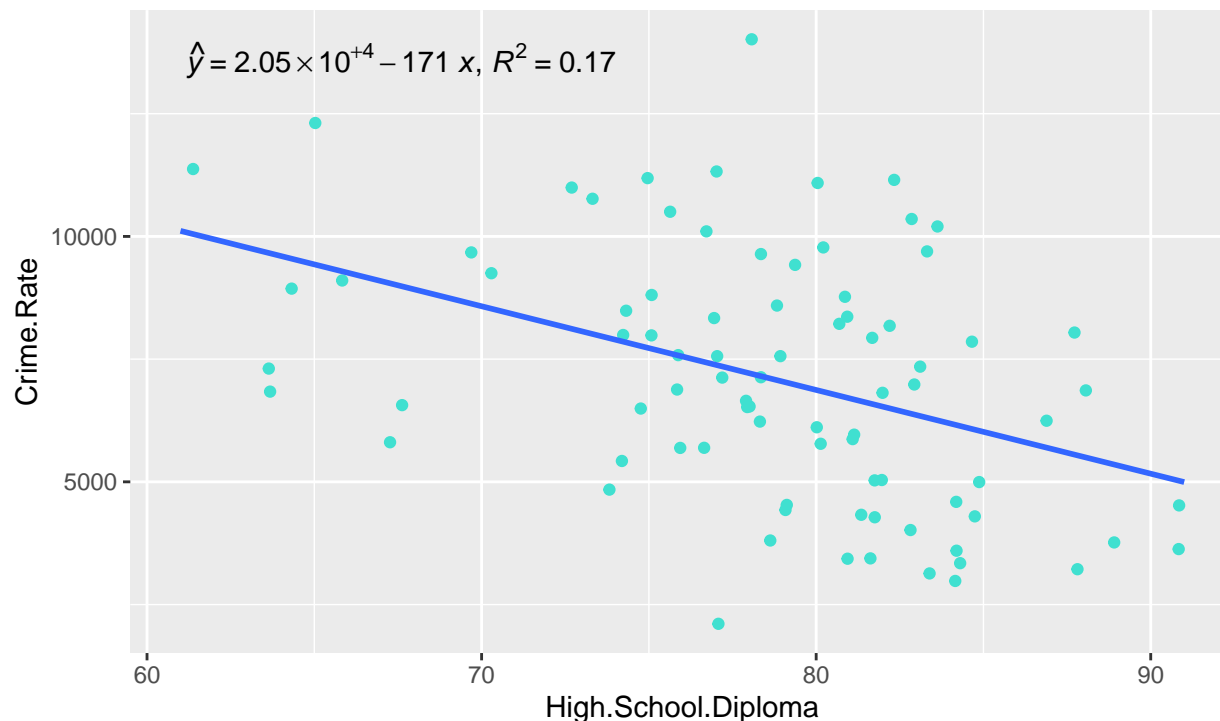
As $p < 0.05$, we reject H_0 at $\alpha = 0.05$ and conclude that there is statistically significant evidence for a linear relationship between crime rate and percent of high school graduates.

h. Give a scatter plot of crime rates vs. percent of high school graduates, with the regression line. Comment about linearity

```
ggplot(educ, aes(x=High.School.Diploma, y=Crime.Rate)) +
  geom_jitter(color="turquoise") +
  geom_smooth(method='lm', formula= y~x,
             se=FALSE,
             show.legend=TRUE) +
  stat_poly_eq(eq.with.lhs = "italic(hat(y))~`=~'",
             use_label(c("eq", "R2")))+
  labs(title = paste("Scatterplot of Crime Rate and Percent of High School Graduates \n with Linear Reg",
                    subtitle = "by Jerry Yu")) +
  theme(plot.title = element_text(size = 12))
```

Scatterplot of Crime Rate and Percent of High School Graduates
with Linear Regression Line and Equation

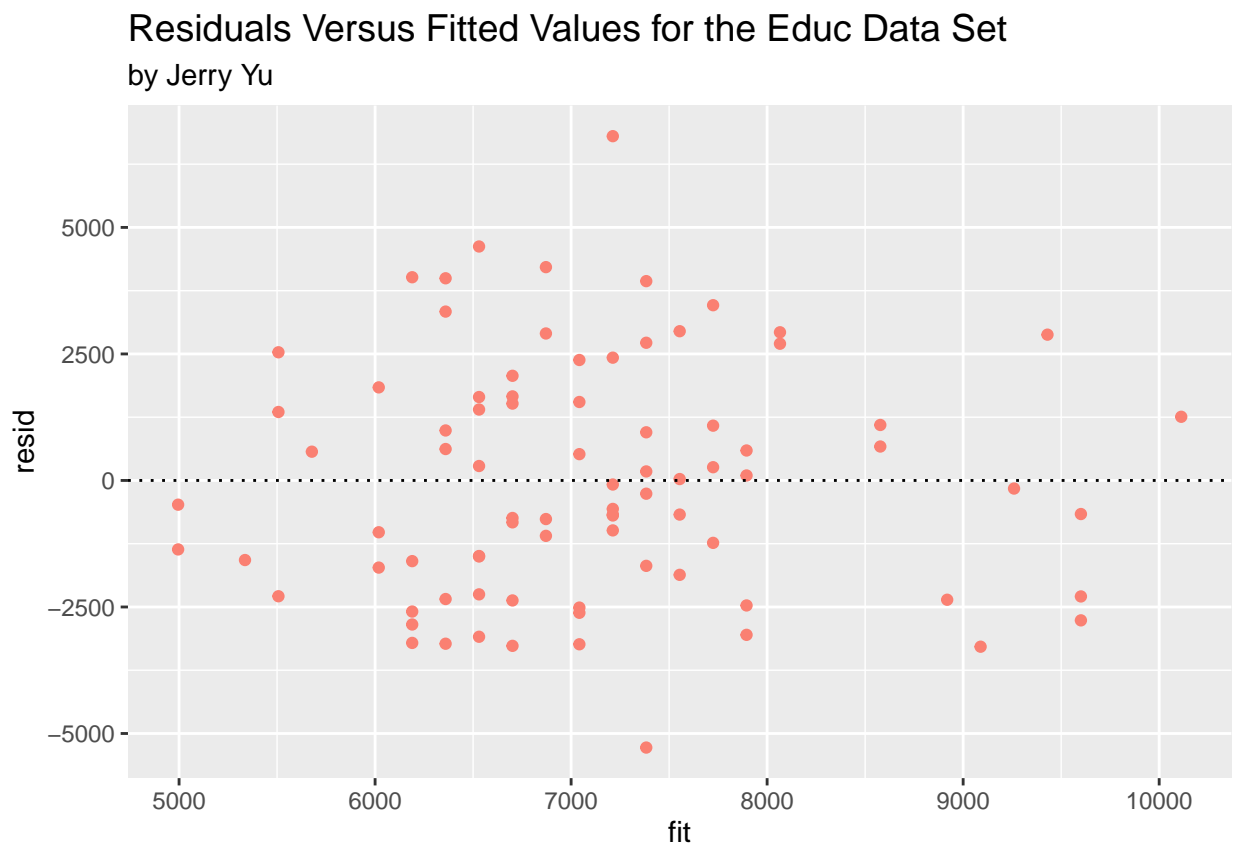
by Jerry Yu



As there does not seem to be a nonlinear pattern in the scatterplot, and the regression line seems to slice across the residuals equally, leaving about 1/2 above and below. I would say that the data seems linear.

i. Give the Residual Plot (residuals vs. fitted values). Test for Non-Linear and Non-constant variance.

```
educmr <- tibble(  
  "fit" = educm$fitted.values,  
  "resid" = educm$residuals  
)  
  
ggplot(educmr, aes(x=fit, y=resid)) +  
  geom_jitter(color="salmon") +  
  geom_hline(yintercept = 0, linetype="dotted") +  
  labs(title = paste("Residuals Versus Fitted Values for the Educ Data Set"),  
       subtitle = "by Jerry Yu") +  
  theme(plot.title = element_text(size = 14))
```



As there do not seem to be patterns in the distribution of the residuals, nor any fan and funnel shapes, I conclude that the variance is likely linear and constant.

j. Conduct Breusch-Pagan Test for the constancy of the error variance. Be sure to give an appropriate null and alternate hypothesis, test statistic, its associated degrees of freedom, and the p-value.

```
educmbp <- bptest(educm, studentize = FALSE)
ncvTest(educm)
```

```
## Non-constant Variance Score Test
## Variance formula: ~ fitted.values
## Chisquare = 0.005045022, Df = 1, p = 0.94338
```

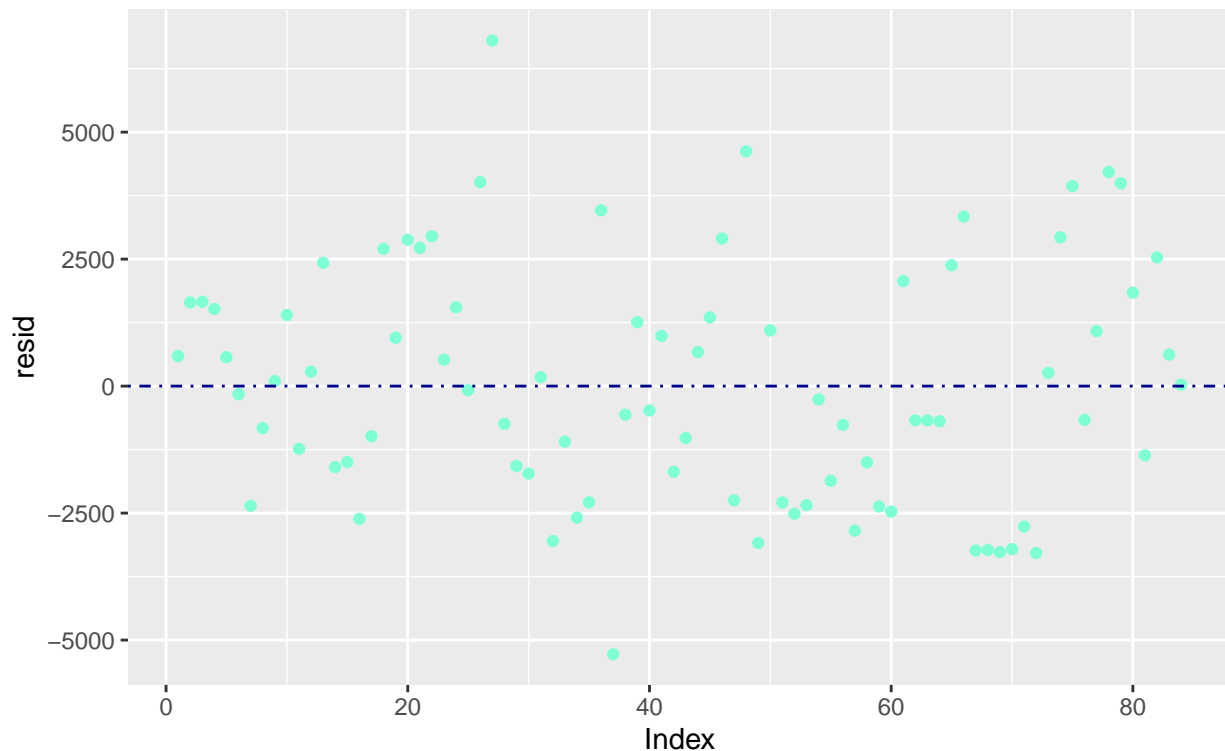
- H0: Equal Variance Among Errors
- HA: Unequal Variance Among Errors
- Degree of Freedom: 1
- P Value: 0.9433752

k. Index Plot to test for Independence of errors.

```
ggplot(educmr, aes(x = 1:length(resid), y = resid)) +
  geom_point(color = "aquamarine") +
  labs(x = "Index",
       title = "Residual Time Sequence Plot for Educ Data",
       subtitle = "by Jerry Yu") +
  geom_hline(yintercept = 0,
            color = "darkblue",
            linetype = "dotdash")
```

Residual Time Sequence Plot for Educ Data

by Jerry Yu



1. Conduct Durbin-Watson Test. Be sure to give an appropriate null and alternate hypothesis, test statistic and the p-value.

```
dwtest(educm)
```

```
##  
## Durbin-Watson test  
##  
## data: educm  
## DW = 1.4951, p-value = 0.008696  
## alternative hypothesis: true autocorrelation is greater than 0
```

```
educmw <-durbinWatsonTest(educm)  
educmw
```

```
## lag Autocorrelation D-W Statistic p-value  
## 1 0.25204 1.495148 0.012  
## Alternative hypothesis: rho != 0
```

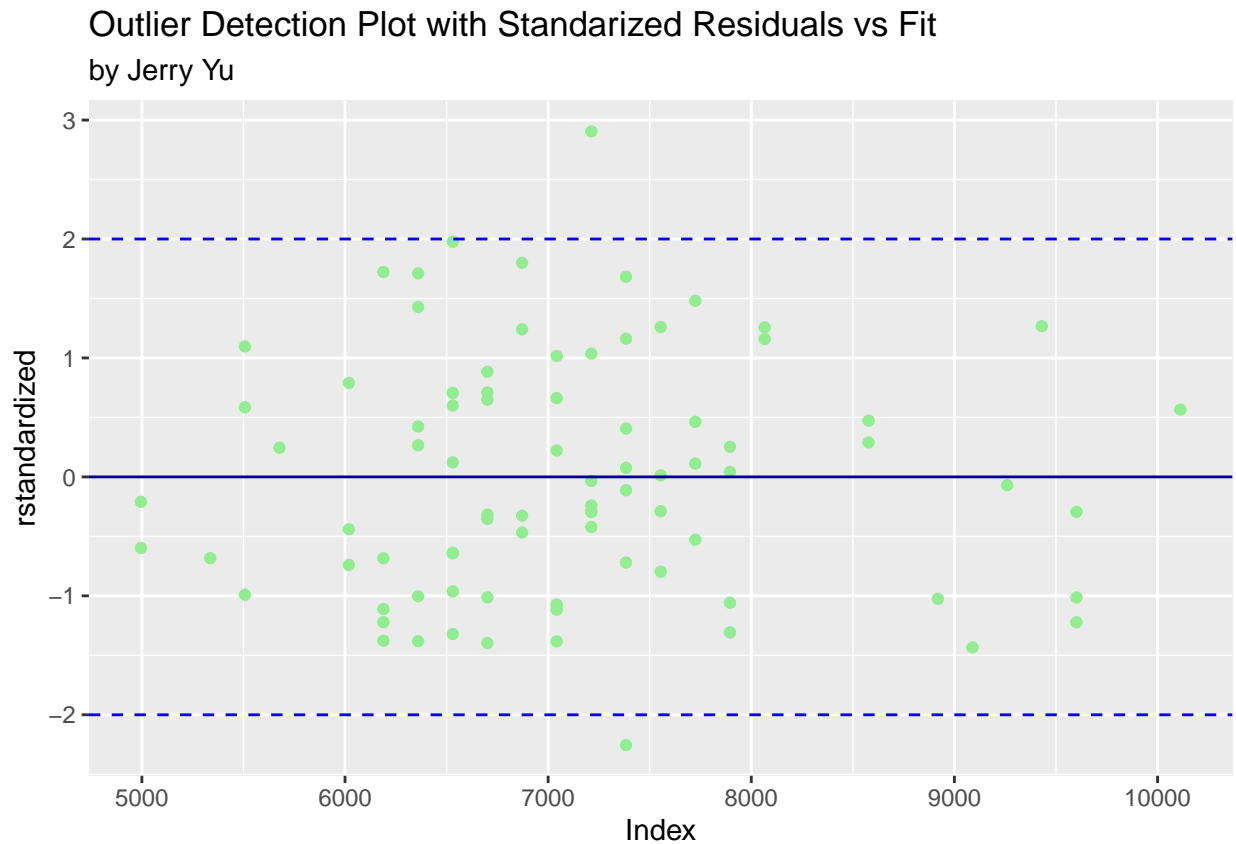
- H0: Errors are uncorrelated over time
- HA: Errors are correlated (either positive or negative). I used the `car` test where the alternative hypothesis is 2 sided

- Test Statistic: 1.4951485
- p value: 0.012

m. Outlier deduction test [Plot standardized Residuals versus fitted values]

```
educmrs <- add_column(educmr, "rstandardized"=rstandard(educm))

ggplot(educmrs, aes(x = fit, y = rstandardized)) +
  geom_point(color = "lightgreen") +
  labs(x = "Index",
       title = "Outlier Detection Plot with Standarized Residuals vs Fit",
       subtitle = "by Jerry Yu") +
  geom_hline(yintercept = 0,
            color = "darkblue",
            linetype = "solid") +
  geom_hline(yintercept = -2,
            color = "blue",
            linetype = "dashed")+
  geom_hline(yintercept = 2,
            color = "blue",
            linetype = "dashed")
```




```
educmro <- filter(educmrs,abs(rstandardized) >2)
educmro
```

```
## # A tibble: 2 x 3
##   fit   resid rstandardized
##   <dbl> <dbl>         <dbl>
## 1 7213.  6803.           2.90
## 2 7383. -5278.          -2.25
```

We have 2 outliers, one where the fit = 7212.7352313 and 7383.31042

n. Give a Histogram of the residuals and the density curve. Comment about the distribution of residuals.

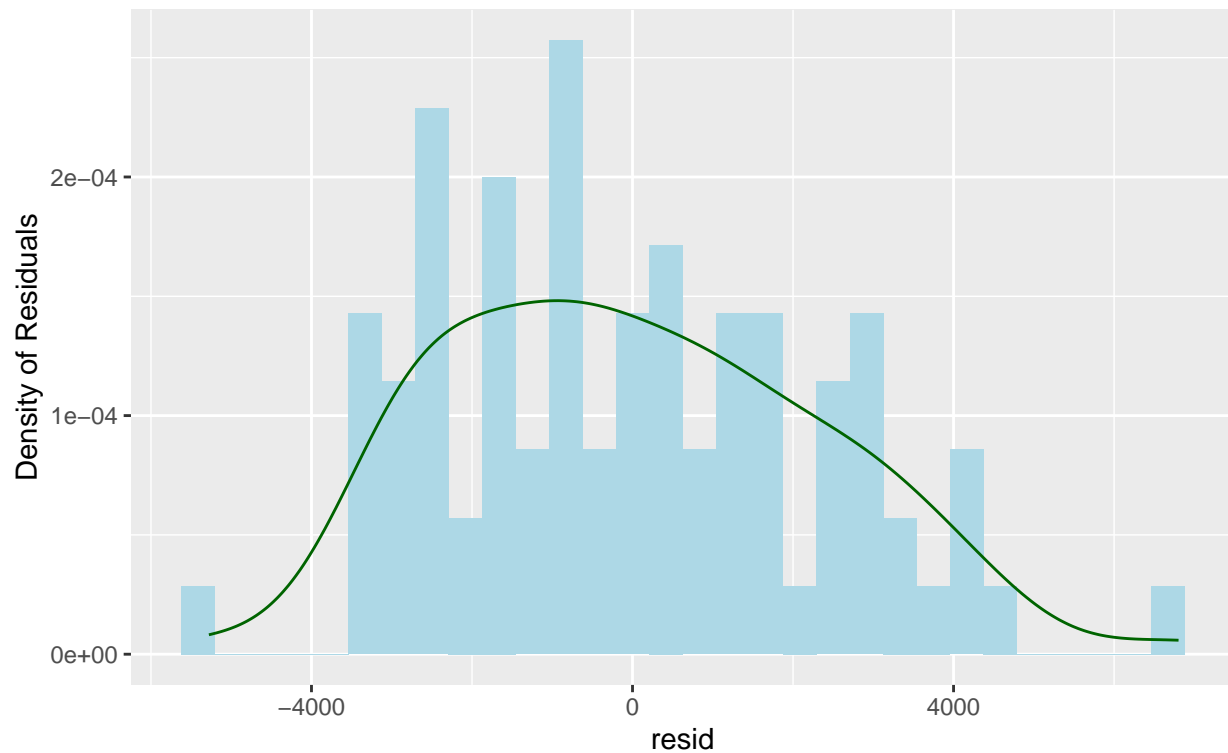
```
ggplot(data = educmrs, aes(x = resid, y = ..density..)) +
  geom_histogram(fill = "lightblue") +
  geom_density(color = "darkgreen") +
  labs(
    title = paste("Histogram and Density Plot of Residuals for Data Set Educ"),
    subtitle = "by Jerry Yu"
  ) +
  ylab("Density of Residuals")
```

```
## Warning: The dot-dot notation (`..density..`) was deprecated in ggplot2 3.4.0.
## i Please use `after_stat(density)` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Histogram and Density Plot of Residuals for Data Set Educ

by Jerry Yu



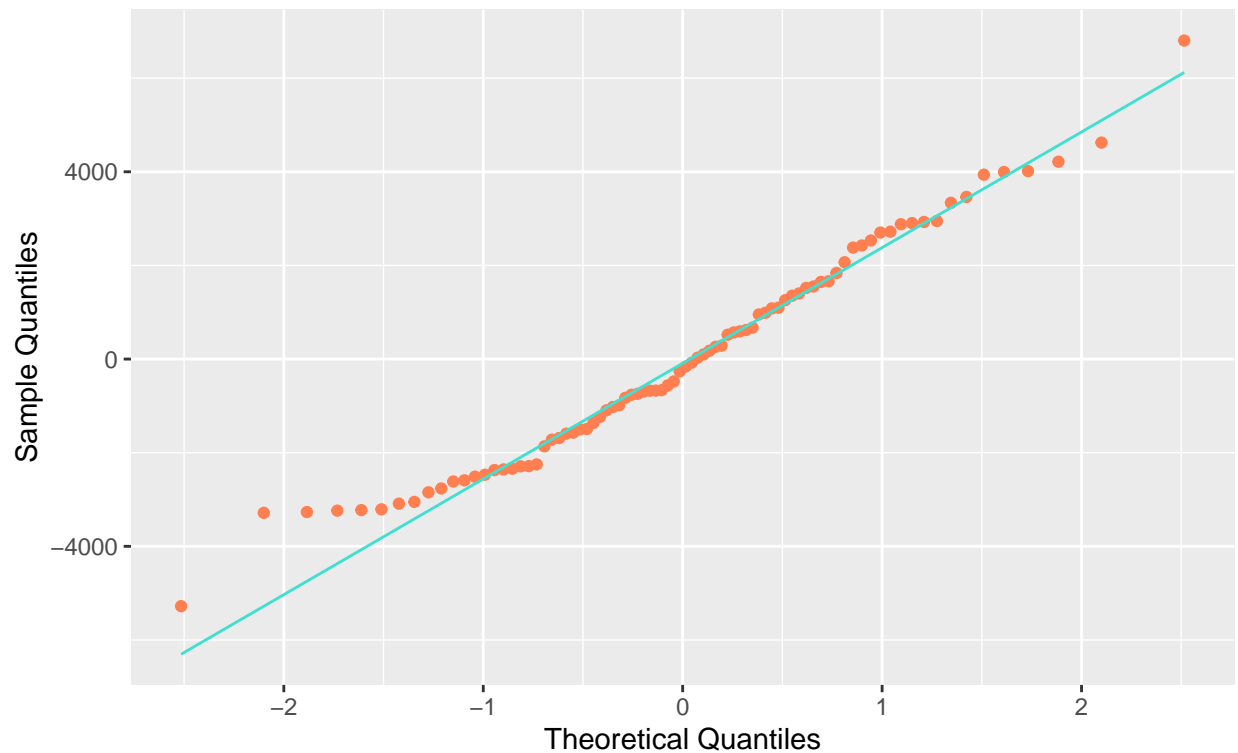
There seems to be a slight right skew in the data, The 2 outliers detected in part m are clearly visible.

o. Give a QQ-plot of the residuals to test for normality of error terms. Comment about the distribution of residuals.

```
ggplot(data = educmrs, aes(sample = resid))+  
  geom_qq( color="coral")+  
  geom_qq_line( color="turquoise")+  
  labs(  
    title = paste("QQ Plot of Residuals for Educ Linear Regression Model"),  
    subtitle = "by Jerry Yu"  
  ) +  
  xlab("Theoretical Quantiles")+  
  ylab("Sample Quantiles")
```

QQ Plot of Residuals for Educ Linear Regression Model

by Jerry Yu



The data visually does not look normal, as the extreme Residuals both look flatter than the Theoretical Residuals (the line).

p. Conduct a Shapiro-Wilk Test on the residuals. Be sure to give an appropriate null and alternate hypothesis, test statistic and the p-value. Give the p-value for this test and explain what this means in terms of our model assumptions.

```
shap1 <- shapiro.test(educmrs$resid)
shap1
```

```
##
##  Shapiro-Wilk normality test
##
## data:  educmrs$resid
## W = 0.97763, p-value = 0.1515
```

- H_0 : The random error is normally distributed
- H_a : The random error is not normally distributed
- Test Statistic: 0.9776328
- p value: 0.1514916

As $p > 0.05$, we fail to reject H_0 at $\alpha = 0.05$ and conclude that there is no statistically significant evidence that the random error is not normally distributed.

Question 2

- a. Give a scatter plot
- b. Find the least squares regression.
- c. Give the Residual Plot (residuals vs. fitted values). Test for Non-Linear and Non-constant variance.
- d. Conduct Breusch-Pagan Test for the constancy of the error variance.
- e. Index Plot to test for Independence of errors.
- f. Conduct Durbin-Watson Test.
- g. outlier deduction test. [Plot standardized Residuals versus fitted values]
- h. Give a Histogram of the residuals.
- i. Give a QQ-plot of the residuals. Normality of error terms.
- j. Conduct a Shapiro-Wilk Test on the residuals. Give the p-value for this test and explain what this means in terms of our model assumptions.
- k. Give the ANOVA Table for this regression analysis. Based on your ANOVA table, is the linear relationship between X and Y statistically significant? Be sure to give an appropriate test statistic, its associated degrees of freedom, and the p-value.

Question 3

- a. Based on your ANOVA table, is the linear relationship between X and Y statistically significant? Be sure to give an appropriate null and alternate hypothesis, test statistic, its associated degrees of freedom, and the p-value.
- b. Give a scatter plot of clot vs. sur_time, with the regression line. Comment about linearity

- c. Give the Residual Plot (residuals vs. fitted values). Test for Non-Linear and Non-constant variance.
- d. Conduct Breusch-Pagan Test for the constancy of the error variance. Be sure to give an appropriate null and alternate hypothesis, test statistic, its associated degrees of freedom, and the p-value.
- e. Index Plot to test for Independence of errors.
- f. Conduct Durbin-Watson Test. Be sure to give an appropriate null and alternate hypothesis, test statistic and the p-value.
- g. Outlier deduction test [Plot standardized Residuals versus fitted values]
- h. Give a Histogram of the residuals and the density curve. Comment about the distribution of residuals.
- i. Give a QQ-plot of the residuals to test for normality of error terms. Comment about the distribution of residuals.
- j. Conduct a Shapiro-Wilk Test on the residuals. Be sure to give an appropriate null and alternate hypothesis, test statistic and the p-value. Give the p-value for this test and explain what this means in terms of our model assumptions