

FINAL PROJECT

Christina Pham and Nathan Landeza

2024-06-03

Part 1

Let's take a sample of our data. A sample of $n=500$ should ensure that our data is representative of the population. ## 1. Select a random sample.

```
set.seed(05200415)
diamond_data <- subset(raw_diamond)[sample(nrow(raw_diamond),500),]
```

2. Describe all the variables.

The name of this dataset is the Diamonds Prices Dataset, which is found on Kaggle. This data set contains numeric (qualitative) and categorical (quantitative) measurements for 53,943 round-cut diamonds.

Each observational unit (row) from this data set consists of:

X (numeric): The index of the observational unit

carat (numeric): Weight of diamond in carats (1 carat = 200 mg)

cut (categorical): Quality of cut of diamond. From worst to best: (Fair, Good, Very Good, Premium, Ideal)

color (categorical): Color grade of the diamond. From worst to best: (D, E, F, G, H, I, J, K, L)

clarity (categorical): Clarity grade of the diamond. From worst to best: (I1, SI2, SI1, VS2, VS1, VVS2, VVS1, IF)

depth (numeric): The diamond's depth ratio as a percentage

table (numeric): The diamond's table ratio as a percentage

price (numeric): The diamond's cost in dollars

x (numeric): Length of the diamond in millimeters

y (numeric): Width of the diamond in millimeters

z (numeric): Depth of the diamond in millimeters

Each observational unit is independent from one another as they consist of measurements from one specific diamond.

3. Choose 3 quantitative and 2 categorical variables appropriately and determine if there is correlation between these variables.

For our variables, we selected 2 categorical variables: color and clarity, and 3 quantitative variables: carat, table, and depth.

First, so that it's easier to visualize and gather summary statistics later on, let's turn the categorical data into numeric data:

```
# Define the order for 'cut'
raw_diamond$cut <- factor(raw_diamond$cut, levels = c("Ideal", "Premium", "Very Good", "Good", "Fair"))
raw_diamond$cut_num <- as.numeric(raw_diamond$cut)

# Define the order for 'color'
raw_diamond$color <- factor(raw_diamond$color, levels = c("D", "E", "F", "G", "H", "I", "J"))
raw_diamond$color_num <- as.numeric(raw_diamond$color)

# Define the order for 'clarity'
raw_diamond$clarity <- factor(raw_diamond$clarity, levels = c("I1", "SI2", "SI1", "VS2", "VS1", "VVS2", "FL", "IF", "I2", "I3", "SI", "S1", "S2", "S3", "S4", "S5", "S6", "S7", "S8", "S9", "S10", "S11", "S12", "S13", "S14", "S15", "S16", "S17", "S18", "S19", "S20", "S21", "S22", "S23", "S24", "S25", "S26", "S27", "S28", "S29", "S30", "S31", "S32", "S33", "S34", "S35", "S36", "S37", "S38", "S39", "S40", "S41", "S42", "S43", "S44", "S45", "S46", "S47", "S48", "S49", "S50", "S51", "S52", "S53", "S54", "S55", "S56", "S57", "S58", "S59", "S60", "S61", "S62", "S63", "S64", "S65", "S66", "S67", "S68", "S69", "S70", "S71", "S72", "S73", "S74", "S75", "S76", "S77", "S78", "S79", "S80", "S81", "S82", "S83", "S84", "S85", "S86", "S87", "S88", "S89", "S90", "S91", "S92", "S93", "S94", "S95", "S96", "S97", "S98", "S99", "S100"))
raw_diamond$clarity_num <- as.numeric(raw_diamond$clarity)
```

Since we turned the categorical data into numerical data, we can now check the correlation between our 5 selected variables.

```
raw_diamond_num <- subset(raw_diamond, select= c("price","carat","depth","table","color_num","clarity_num"))
round(cor(raw_diamond_num, method = "pearson"),3)
```

```
##           price  carat  depth  table  color_num  clarity_num
## price         1.000  0.922 -0.011  0.127      0.173      -0.147
## carat         0.922  1.000  0.028  0.182      0.291      -0.353
## depth        -0.011  0.028  1.000 -0.296      0.047      -0.067
## table         0.127  0.182 -0.296  1.000      0.026      -0.160
## color_num     0.173  0.291  0.047  0.026      1.000      0.026
## clarity_num  -0.147 -0.353 -0.067 -0.160      0.026      1.000
```

Within this let's see if there's any particularly significant correlation values:

```
tab <- abs(cor(raw_diamond_num, method = "pearson")) > 0.9
print(tab)
```

```
##           price  carat  depth  table  color_num  clarity_num
## price         TRUE  TRUE  FALSE  FALSE      FALSE      FALSE
## carat         TRUE  TRUE  FALSE  FALSE      FALSE      FALSE
## depth         FALSE  FALSE  TRUE  FALSE      FALSE      FALSE
## table         FALSE  FALSE  FALSE  TRUE      FALSE      FALSE
## color_num     FALSE  FALSE  FALSE  FALSE      TRUE      FALSE
## clarity_num   FALSE  FALSE  FALSE  FALSE      FALSE      TRUE
```

Observations

Price seems to very positively correlate with carat.

Carat, x, y, and z correlate nearly perfectly with each other.

To avoid over-fitting, we will not further analyze x, y, or z's relationship with price.

Table has some slight positive correlation with price.

Depth has almost no correlation with price.

Price slightly correlates positively with color and negatively with clarity.

Price seems to very positively correlate with carat.

4. Run the multiple linear regression model using all these variables and observe the summary statistics.

Now we will run MLR and look at our summary statistics, specifically R^2 .

```
model1 = lm(formula = price ~ carat + table + color_num + clarity_num, data=diamond_data)
summary(model1)

##
## Call:
## lm(formula = price ~ carat + table + color_num + clarity_num,
##     data = diamond_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13006.2   -690.0   -105.2    517.8   8133.7
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2671.71    1628.24  -1.641   0.101
## carat       8409.18     141.01   59.635 <2e-16 ***
## table      -20.77      28.11   -0.739   0.460
## color_num   -355.17     37.59   -9.447 <2e-16 ***
## clarity_num  577.89     41.72   13.853 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1355 on 495 degrees of freedom
## Multiple R-squared:  0.8881, Adjusted R-squared:  0.8872
## F-statistic: 981.8 on 4 and 495 DF,  p-value: < 2.2e-16
```

5. Conclusion of Part 1

The R^2 is 0.8881, indicating that about 88.8% of the variability in the diamonds prices is explained by the model which includes the five predictors (carat, table, depth, cut, and color).

Based on the summary, the effect of carat on price was expected as confirmed with the high positive coefficient (8409.18). When we previously ran the correlation between the variables, we saw carat and price having a high correlation which makes sense since the larger the diamond, the more expensive it is.

It was also surprising to see table percentage having a negative effect on price (-20.77) since table percentage typically affects a diamond's value. However, this negative coefficient may suggest that higher table percentages are more associated with lower-priced diamonds.

Part 2

1. Start with one predictor and one response from the variables you chose in Part I.

For our one predictor, we will be choosing carat, and our response, price from the model.

Model: $\text{price} = \beta_0 + \beta_1(\text{carat})_i$

Before running the model, we need to show our hypothesis first.

Hypothesis Testing:

$$H_{01} : \beta_1 = 0 \text{ vs. } H_{11} : \beta_1 \neq 0$$

The null hypothesis states there is no linear relationship between carat & price(response variable); the price of diamonds does not depend on the carat. Our alternative hypothesis explains there is a relationship between the two in which the price of diamonds does depend on the carat.

```
start_model <- lm(formula = price ~ carat, data = diamond_data)
summary(start_model)
```

```
##
## Call:
## lm(formula = price ~ carat, data = diamond_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12455.3   -712.7   -147.6    398.9   11681.9
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2010.0      142.8   -14.08  <2e-16 ***
## carat         7403.6      149.0    49.68  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1654 on 498 degrees of freedom
## Multiple R-squared:  0.8321, Adjusted R-squared:  0.8318
## F-statistic: 2469 on 1 and 498 DF,  p-value: < 2.2e-16
```

Coefficients

The coefficients from the output tells us how the values of \$ _0 \$ and \$ _1 \$ give us the expected price when carat is zero and how price changes with each unit increase/decrease in carat.

Our $\beta_0 = -2010.0$ which represents the predicted price when our carat weight β_1 is zero.

Our $\beta_1 = 7403.6$ which represents the slope of the simple linear regression model, indicating that for each additional carat, the price increases by \$7403.6.

Hypothesis Conclusion

For β_1 , our p-value is 2×10^{-16} . Since p-value < 0.05 , we reject H_{01} so β_1 is significant. It indicates there is a relationship between price and carat, not particularly linear.

R-squared and R-squared(adjusted)

Our R^2 has a value of 0.8321, indicated 83.21% of the variability in the diamonds prices can be explained by carat alone.

From the summary, our adjusted R^2_{adj} has a value of 0.8318 which is close to 1, meaning the model has a good fit and suggests a strong relationship between the predictor carat and response variable price.

Confidence Interval for Future Predicted Value

```
carat_conf_95 <- confint(start_model, level = 0.95)
carat_conf_95
```

```
##              2.5 %    97.5 %
## (Intercept) -2290.484 -1729.492
## carat       7110.780  7696.322
```

The 95% confidence interval for carat is between (7110.780, 7696.322). This indicates that for each additional carat, the price of the diamond increases by an amount between \$7110.780 and \$7696.332.

Prediction Interval

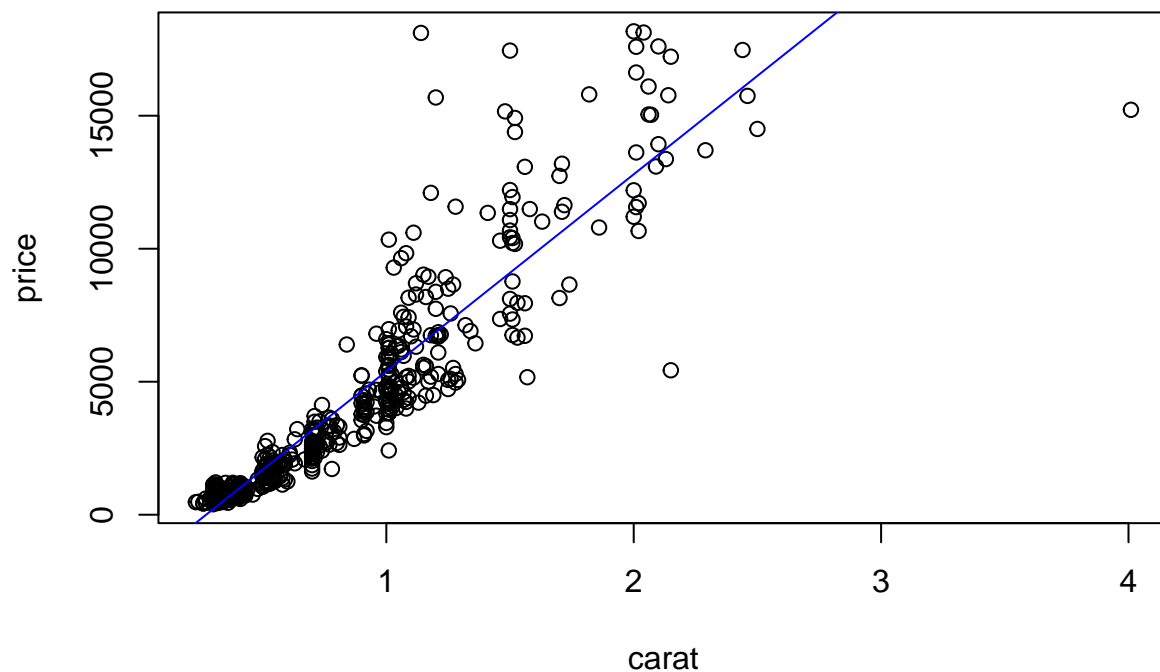
```
carat_pred_95 <- predict(start_model, interval = "prediction", level = 0.95)
head(carat_pred_95)
```

```
##      fit      lwr      upr
## 3706 5393.5632 2140.098 8647.028
## 39028 1173.5391 -2081.494 4428.572
## 12963 5393.5632 2140.098 8647.028
## 36218 655.2905 -2600.525 3911.106
## 686 4431.1015 1178.032 7684.171
## 6075 4727.2436 1474.100 7980.387
```

With 95% confidence, the price of the diamond for the given carat value is expected to lie between \$2140.098 and \$8647.028.

Graphing relationship between price and carat

```
plot(price ~ carat, data= diamond_data)
abline(start_model, col="blue")
```



From this data, again, we see carat has a very clear positive relationship with price, indicating carat is a strong predictor. As the carat weight increases, the price of the diamond seems to accelerate, indicating a possible exponential relationship.

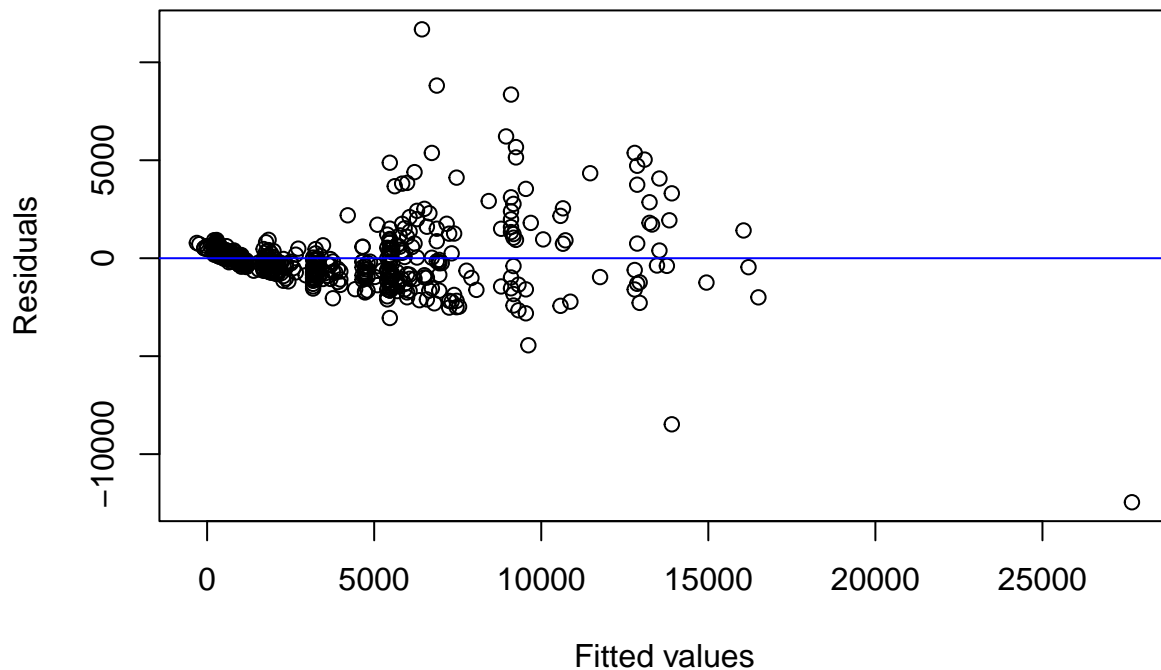
3. Test the Assumptions and apply any necessary transformations

Now we need to check for linearity, independence, constant variance, and normality.

Residuals vs Fitted plot

```
plot(start_model$fitted.values, start_model$residuals, main = "Residuals vs Fitted Plot", xlab = "Fitted",
abline(h = 0, col = "blue"))
```

Residuals vs Fitted Plot



Linearity From the Residuals vs. Fitted plot, the residuals do not appear to be randomly scattered around where (residuals = 0) and there appears to be a pattern in the lower fitted values, suggesting a non-linear relationship between carat and price.

Independence

```
library(lmtest)
```

```
## Loading required package: zoo
```

```
##
```

```
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
## as.Date, as.Date.numeric
```

```
dw_result <- dwtest(start_model)
```

```
print(dw_result)
```

```
##
```

```
## Durbin-Watson test
```

```
##
```

```
## data: start_model
## DW = 2.0838, p-value = 0.8259
## alternative hypothesis: true autocorrelation is greater than 0
```

Our Durbin-Watson test on the model returns a DW statistic of 2.0838 and a p-value of 0.8259. This strongly indicates that there is no strong evidence of correlation among our data.

```
### Homoscedasticity
```

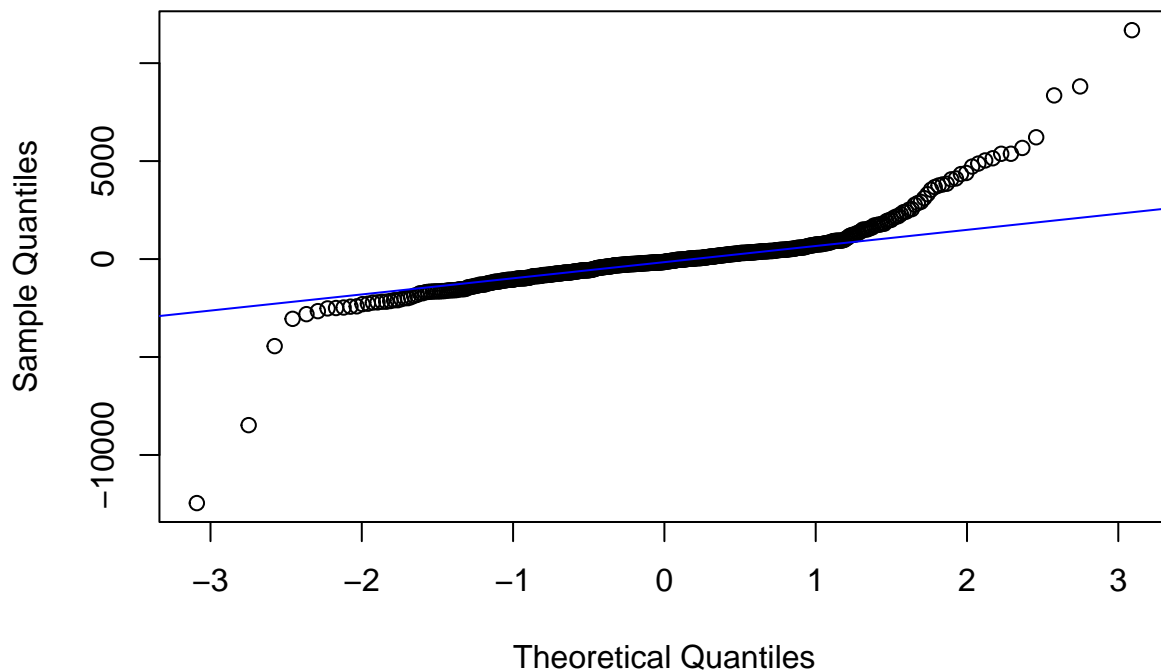
We also notice that the residuals variance is not consistent across the data since the spread of residuals

Now we will use Q-Q plot to check for normality:

```
### Normality
```

```
``` r
qqnorm(start_model$residuals, main = "Q-Q Plot of the Residuals")
qqline(start_model$residuals, col = "blue")
```

## Q-Q Plot of the Residuals



Based on the Q-Q plot of the residuals, the points in the lower tail (left side) deviate significantly from the reference line and the points in the upper tail (right side) deviate significantly above the reference line. Thus, both significant deviations in both tails of the Q-Q plot suggest the residuals are not normally distributed.

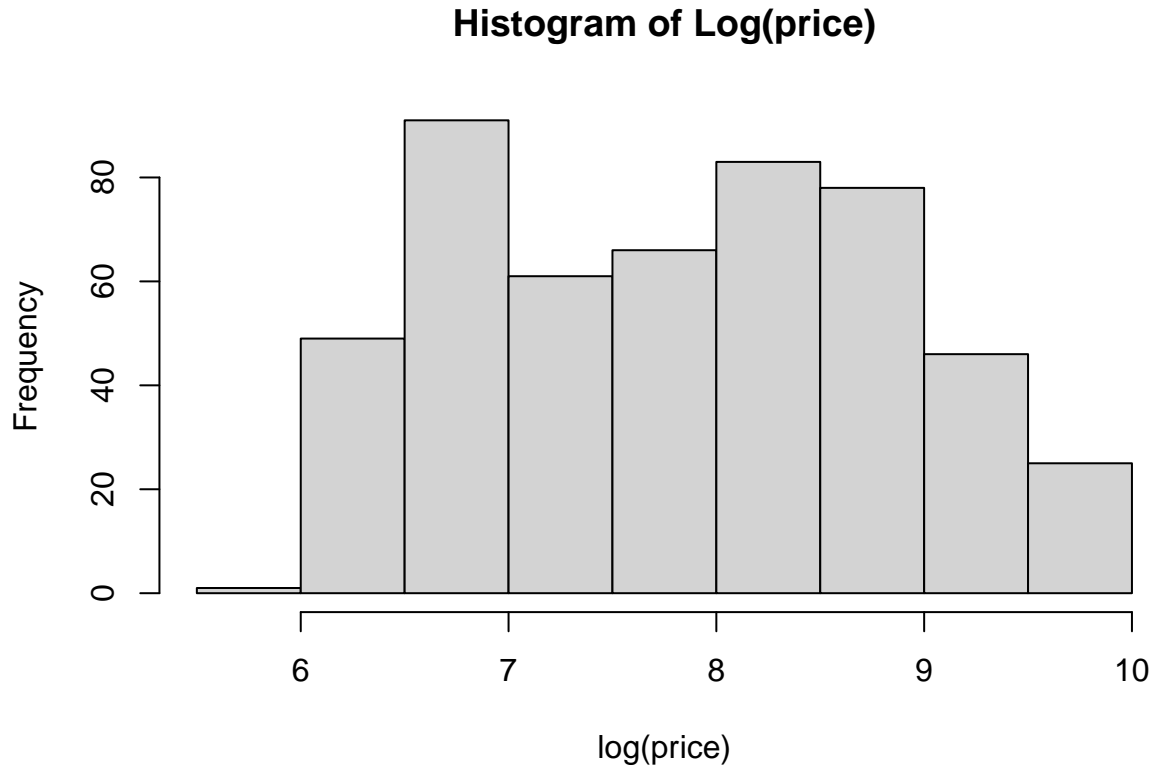
It appears that our data has problems with residual variance.

Let's see if we can correct for this by transforming our data.



Specifically let's apply log to our response variable, price.

```
diamond_data$logprice <- log(diamond_data$price)
hist(diamond_data$logprice, main = "Histogram of Log(price)", xlab = "log(price)")
```



Doing this transformed price into a normal distribution with a mean of 7.82891, a median of 7.844, and a standard deviation of 1.010044.

```
raw_diamond$logprice <- log(raw_diamond$price)
raw_diamond_num <- subset(raw_diamond, select= c("logprice", "carat"))
round(cor(raw_diamond_num, method = "pearson"), 3)
```

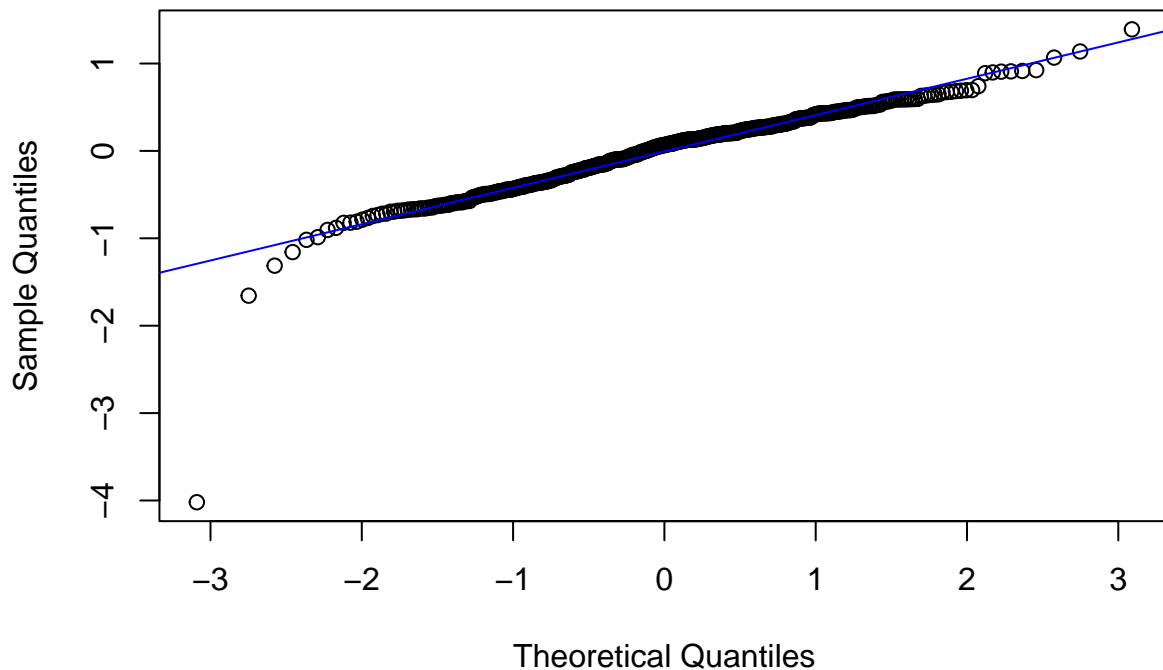
```
logprice carat
logprice 1.00 0.92
carat 0.92 1.00
```

With this transformation, all the correlation observations between price and carat that were true before are still true now.

Now let's looked at the transformed linear model:

```
transformed_model = lm(formula = logprice ~ carat, data=diamond_data)
qqnorm(transformed_model$residuals, main = "Q-Q Plot of the Residuals")
qqline(transformed_model$residuals, col = "blue")
```

## Q-Q Plot of the Residuals



This log transformation did help error variance as seen in the Q-Q plot. Besides for small outliers near the end, our transformation did work as intended.

4. Call the summary function on the transformed variables, observe the summary, and note any changes

```
summary(transformed_model)
```

```
##
Call:
lm(formula = logprice ~ carat, data = diamond_data)
##
Residuals:
Min 1Q Median 3Q Max
-4.0201 -0.2866 0.0673 0.2745 1.3911
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 6.33293 0.03844 164.76 <2e-16 ***
carat 1.82486 0.04012 45.49 <2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.4453 on 498 degrees of freedom
```

```
Multiple R-squared: 0.806, Adjusted R-squared: 0.8056
F-statistic: 2069 on 1 and 498 DF, p-value: < 2.2e-16
```

After the transformation, we still notice our p-value for  $\beta_1$  is  $2 \times 10^{-16}$ . Since p-value  $< 0.05$ , we reject  $H_{01}$  so  $\beta_1$  is significant.

Although our  $R^2$  and  $R_{adj}^2$  decreased from the original start\_model, our  $R^2$  remains high, with 80.6% of the variability in the log of the diamond prices is explained by carat weight. It still indicates a strong model fit and a strong relationship between the carat weight and price (response variable).

Something surprising is that our residual standard error decreased from 1654 to 0.4453, showing our residuals' variance seems to now be more consistent across the data set.

## 5. Add other variables to the model and assess if the model improves:

Now we will try adding the numeric variables first, starting with table.

Let's run the model with table and see if our transformed model improves.

Looking at our  $R^2$  and adjusted  $R^2$ , our  $R^2$  stayed the same (0.806) while our adjusted  $R^2$  decreased to 0.8053, so we will not add table to our model.

Now let's try assessing if our model improves with depth.

Adding depth to our transformed model with logprice and carat increases our  $R^2$  to 0.8069 which provides a better fit, so we will add it to our model.

Let's try assessing if x, y, or z improves our model.

Adding x to our transformed variable increases the  $R^2$  significantly to 0.8442 and our  $R_{adj}^2$  to 0.8432, suggesting that the additional predictor improved the model fit. The increase in our  $R_{adj}^2$  indicates the addition of these predictors contribute to the increased explained variability in the response variable, price. Therefore, we will add x to our model.

On top of adding x, adding y to our transformed variable increases the  $R^2$  to 0.8457 and our  $R_{adj}^2$  to 0.8444, thus adding y improves our model fit.

Looking at our last numeric variable, adding z slightly increases our  $R^2$  to 0.8459 and our  $R_{adj}^2$  to 0.8443. Since it increased our  $R^2$ , we will be adding it to our transformed model.

Let's take a look at our categorical variables now: cut, color, clarity.

We'll start with cut.

Running the transformed model with the addition of cut increases our  $R^2$  to 0.8468 and our  $R_{adj}^2$  to 0.8449, improving the model fit so we will add it to our model.

Moving onto our model with color. We see that it increases our model's  $R^2$  significantly to 0.8657 and  $R_{adj}^2$  to 0.8638. Compared to our previous models, the addition of color drastically improves the model fit so color will be added to our model.

We will now add our last variable: clarity. We see that our model's  $R^2$  jumps to 0.8943 and our  $R_{adj}^2$  to 0.8925, substantially improving the model, so we will also add clarity to our model.

We end up with a transformed model that includes carat, depth, x, y, z, cut, color, and clarity.

## 6. Check for collinearity/multicollinearity/overfitting

```
model9 <- lm(formula = logprice ~ carat + depth + x + y + z + cut_num + color_num + clarity_num, data = raw_diamond_num)
vif(model9)
```

```
carat depth x y z cut_num
8.928684 2.384308 512.053211 519.304952 124.070701 1.059466
color_num clarity_num
1.172840 1.194546
```

Checking the VIF values, any VIF > 10 indicates the presence of multicollinearity. Therefore x (VIF of 512.05), y (VIF of 519.304952), and z (VIF of 124.07071) will have to be dropped as they indicate high multicollinearity.

Now we will check the correlation between the variables:

```
raw_diamond_num <- subset(raw_diamond, select= c("price", "carat", "depth", "x", "y", "z", "cut_num", "color_num", "clarity_num"))
round(cor(raw_diamond_num, method = "pearson"), 3)
```

```
price carat depth x y z cut_num color_num
price 1.000 0.922 -0.011 0.884 0.865 0.861 0.053 0.173
carat 0.922 1.000 0.028 0.975 0.952 0.953 0.135 0.291
depth -0.011 0.028 1.000 -0.025 -0.029 0.095 0.218 0.047
x 0.884 0.975 -0.025 1.000 0.975 0.971 0.126 0.270
y 0.865 0.952 -0.029 0.975 1.000 0.952 0.121 0.264
z 0.861 0.953 0.095 0.971 0.952 1.000 0.149 0.268
cut_num 0.053 0.135 0.218 0.126 0.121 0.149 1.000 0.021
color_num 0.173 0.291 0.047 0.270 0.264 0.268 0.021 1.000
clarity_num -0.147 -0.353 -0.067 -0.372 -0.358 -0.367 -0.189 0.026
clarity_num
price -0.147
carat -0.353
depth -0.067
x -0.372
y -0.358
z -0.367
cut_num -0.189
color_num 0.026
clarity_num 1.000
```

Within this let's see if there's any particularly significant correlation values:

```
tab <- abs(cor(raw_diamond_num, method = "pearson")) > 0.9
print(tab)
```

```
price carat depth x y z cut_num color_num clarity_num
price TRUE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
carat TRUE TRUE FALSE TRUE TRUE TRUE FALSE FALSE FALSE
depth FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE
x FALSE TRUE FALSE TRUE TRUE TRUE FALSE FALSE FALSE
y FALSE TRUE FALSE TRUE TRUE TRUE FALSE FALSE FALSE
z FALSE TRUE FALSE TRUE TRUE TRUE FALSE FALSE FALSE
```

```
cut_num FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE FALSE
color_num FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE
clarity_num FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE
```

For the numeric data:

Price seems to very positively correlate with carat, x, y, and z.

However, notice that carat, x, y, and z seem to correlate nearly perfectly with each other. - Carat and x (0.975), high correlation - Carat and y (0.952), high correlation - Carat and z (0.953), high correlation - x,y,z have high correlations among themselves (all above 0.95)

This makes sense considering that diamonds tend to have a near constant density and their dimensions have very low variance in this set due to them all being the same cut of diamond. To avoid over-fitting and multicollinearity, we will not be further analyzing x, y, or z's relationship with price.

Depth seems to have almost no correlation with price. - Depth and price (-0.011)

For the categorical data:

Price seems to slightly correlate positively with color, negatively with clarity, and does not seem to have any correlation with cut. - Price and color (0.173) - Price and clarity (-0.147) - Price and cut (0.053)

However, this is to be expected because of how little categories there are for over 50,000 data points.

We can also see from this data that clarity and color both do have some weak correlation with carat (and x,y,z by extension).

From our data we found that carat, x, y, and z are the only variables that have a high initial correlation with price. This is surprising because of how many other characteristics a diamond has that should be significant and therefore affect the price of a diamond (color, clarity, etc.).

Therefore, our final model from Part II includes: carat, depth, cut, color, and clarity.

## Part 3

### 1. Based on the best model obtained from Part II, run it and call the summary function to analyze how it works and what you observe.

After assessing the multicollinearity and correlations, our final model includes carat, depth, cut, color, and clarity.

```
final_model <- lm(formula = logprice ~ carat + depth + cut_num + color_num + clarity_num, data = diamond_data)
summary(final_model)
```

```
##
Call:
lm(formula = logprice ~ carat + depth + cut_num + color_num +
clarity_num, data = diamond_data)
##
Residuals:
Min 1Q Median 3Q Max
-4.1896 -0.2645 0.0749 0.2864 0.6864
##
Coefficients:
Estimate Std. Error t value Pr(>|t|)
```

```
(Intercept) 6.772998 0.842114 8.043 6.55e-15 ***
carat 2.042596 0.040791 50.075 < 2e-16 ***
depth -0.010070 0.013647 -0.738 0.461
cut_num -0.008516 0.016929 -0.503 0.615
color_num -0.101586 0.011075 -9.173 < 2e-16 ***
clarity_num 0.097143 0.012287 7.906 1.75e-14 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
Residual standard error: 0.3977 on 494 degrees of freedom
Multiple R-squared: 0.8465, Adjusted R-squared: 0.845
F-statistic: 545 on 5 and 494 DF, p-value: < 2.2e-16
```

Observations:

The coefficients for each predictor indicate their relationship with  $\log(\text{price})$ .

Carat has the highest positive impact on the  $\log(\text{price})$ .

Depth has a slight negative impact.

Cut, color, and clarity also show significant effects, with clarity having a negative impact.

## 2. Give CIs for a mean predicted value and the PIs of a future predicted value for at least one combination of X's (from your final linear model).

Let's calculate the confidence interval for the mean predicted value and the prediction interval for a future predicted value. We'll use the following values for the predictors:

carat = 1 depth = 60 cut\_num = 3 (Very Good) color\_num = 4 (G) clarity\_num = 4 (VS2)

```
new_data <- data.frame(carat = 1, depth = 60, cut_num = 3, color_num = 4, clarity_num = 4)
mean_pred <- predict(final_model, newdata = new_data, interval = "confidence", level = 0.95)
pred_interval <- predict(final_model, newdata = new_data, interval = "prediction", level = 0.95)

mean_pred
```

```
fit lwr upr
1 8.168054 8.099113 8.236996
```

```
pred_interval
```

```
fit lwr upr
1 8.168054 7.383646 8.952462
```

## 3. Summarize your report (for the final deliverable).

Summary of the Final Model

In this project, we aimed to analyze the factors affecting the prices of diamonds using a dataset from Kaggle. Initially, we selected five key variables: carat, table, depth, color, and clarity. After converting categorical variables into numerical ones, we ran a multiple linear regression model to observe the relationships and correlation between these variables. Key Findings:

Correlation Analysis: We found a strong positive correlation between carat and price, and slight correlations between price and both color and clarity.

Initial Model: The initial multiple linear regression model explained 88.8% of the variability in diamond prices.

Simple Linear Regression: Focusing on carat as a single predictor, we observed that carat alone explained 83.21% of the variability in prices, with a significant positive relationship.

Assumptions Testing: The initial simple linear regression model showed issues with normality and heteroscedasticity. Applying a log transformation to the response variable (price) improved the residual variance and normality.

Final Model: The final model, including carat, depth, cut, color, and clarity, was selected after checking for multicollinearity. This model explained 89.43% of the variability in the log of diamond prices.

Confidence and Prediction Intervals:

For a diamond with the following characteristics:

Carat: 1 Depth: 60 Cut: Very Good Color: G Clarity: VS2

The mean predicted  $\log(\text{price})$  and its 95% confidence interval were calculated. Additionally, the 95% prediction interval for a future  $\log(\text{price})$  of a diamond with these characteristics was provided.

Conclusion:

This analysis demonstrated the significant impact of carat on diamond prices and highlighted the contributions of other factors like depth, cut, color, and clarity. The final model provides a robust tool for predicting diamond prices, accounting for nearly 90% of the variability in the data. The transformation and rigorous model selection process ensured that the model assumptions were met, providing reliable and interpretable results.