



"Analysis and Prediction of Diamond Prices using Regression"

Course: "DADS6001 Applied Modern Statistical Analysis"

Student Name: "Miss Nutthida Yotaprasert"

Student ID: "6610422019"

Supervisor: "Asst. Prof. Dr. Ramidha Srihera"

Project Description: "The project focuses on analyzing and predicting diamond prices by utilizing a diverse dataset to create a statistical model that can effectively forecast diamond values based on key characteristics."

Date: "9 May 2024"

Analysis and Prediction of Diamond Prices using Regression

Introduction

Analyzing data to predict dependent variables from independent variables is crucial in understanding and forecasting data behavior across various fields. Predicting diamond prices based on their characteristics using regression is a compelling example. By utilizing regression, statistical models can be created to predict diamond prices from their various characteristics.

This research project focuses on analyzing diamond price data with a diverse set of characteristics like weight, cut quality, color, clarity, and size. The primary goal is to develop a suitable statistical model to predict diamond prices accurately, aiding analysts and stakeholders in better understanding and predicting diamond values.

The dataset used in the project comprises crucial information such as diamond prices in USD, carat weight, cut quality, color, clarity, and dimensions, essential for effective analysis and prediction of diamond prices. Through data analysis and the creation of appropriate statistical models, a better understanding and accurate prediction of diamond values can be achieved.

Therefore, analyzing data to predict diamond prices showcases the importance and benefits of regression in solving mathematical problems in daily life and predicting crucial variables for decision-making in the diamond industry today. The aim is for the results of this project to be beneficial and provide guidance for future research and development efforts.

Objectives

This project aims to analyze and predict diamond prices using a dataset with diverse characteristics like weight, cut quality, color, clarity, and size. The objective is to create a suitable statistical model to predict diamond prices accurately from these characteristics. By analyzing data and developing appropriate statistical models, a more efficient and precise prediction of diamond values can be achieved. Analyzing and predicting diamond prices demonstrates the importance and benefits of regression in solving mathematical problems in daily life and predicting essential variables for decision-making in the current diamond industry.

References

<https://www.kaggle.com/datasets/shivam2503/diamonds>

https://github.com/Jeniejean/Applied-Stat/blob/main/Diamond_Data.csv

The dataset consists of the following variables:

- Price: The price of the diamond in United States dollars (USD), ranging from \$326 to \$18,823.
- Carat: The weight of the diamond in carats, ranging from 0.2 to 5.01.
- Cut: The quality of the diamond cut (Fair, Good, Very Good, Premium, Ideal), with corresponding numerical values: "Fair" = 0, "Good" = 1, "Very Good" = 2, "Premium" = 3, and "Ideal" = 4.
- Color: The color grade of the diamond, ranging from J (worst) to D (best), with numerical values assigned as follows: J = 0, I = 1, H = 2, G = 3, F = 4, E = 5, D = 6.
- Clarity: The clarity grade of the diamond (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best)), with numerical values assigned as follows: I1 = 0, SI2 = 1, SI1 = 2, VS2 = 3, VS1 = 4, VVS2 = 5, VVS1 = 6, IF = 7.
- x: The length of the diamond in millimeters, ranging from 0 to 10.74.
- y: The width of the diamond in millimeters, ranging from 0 to 58.9.
- z: The depth of the diamond in millimeters, ranging from 0 to 31.8.
- Depth: The total depth percentage of the diamond, calculated as $z / \text{mean}(x, y)$ or $2 * z / (x + y)$, with values ranging from 43 to 79.
- Table: The width of the top of the diamond compared to its widest point, with values ranging from 43 to 95.

R Programming

1. Plotting data and initiating data analysis

- Using the readr package to load and read data from a CSV file from a specified URL
- Using the pairs() function to create a scatter plot matrix to analyze the data distribution in the created data frame

```
library(readr)
data <- read_csv("https://raw.githubusercontent.com/Jeniejean/Applied-Stat/main/Diamond_Data%20(3).csv")
dataf <- data.frame(data$price, data$carat, data$cut, data$color, data$clarity, data$depth, data$table, data$x, data$y, data$z)
pairs(dataf)
```

Rows: 203 Columns: 10

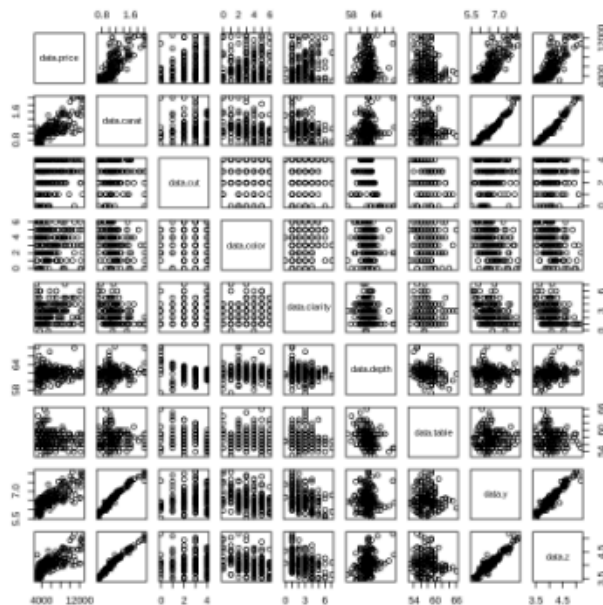
— Column specification —————

Delimiter: ",",

dbl (10): carat, cut, color, clarity, depth, table, price, x, y, z

i Use `spec()` to retrieve the full column specification for this data.

i Specify the column types or set `show_col_types = FALSE` to quiet this message.



2. Variable Selection

Selecting important variables to create the best Linear Regression model from the previously built Linear Regression model

Equation: The linear regression equation derived from the final model is:

$$\text{Diamond Price} = -50330.66 + 121.63 * \text{cut} + 477.39 * \text{color} + 777.89 * \text{clarity} + 304.82 * \text{depth} + 5106.73 * y$$

Adjusted R-squared: The adjusted R-squared value of this model is 0.8756, indicating that the 5 independent variables (cut, color, clarity, depth, y) can explain approximately 87.56% of the variance in diamond prices.

P-value: The p-value of this model is $< 2.2e-16$, which is less than the standard statistical significance value set at 0.05.

Hypotheses: $H_0: \beta_i = 0$
 $H_1: \beta_i \neq 0$

Results Analysis: With a p-value less than 0.05, we can reject the null hypothesis H_0 and accept the alternative hypothesis H_1 , indicating a response of diamond prices to the input independent variables.

Conclusion: The analysis reveals that the 5 independent variables (cut, color, clarity, depth, y) significantly influence diamond prices. The linear regression model developed has a strong ability to explain the data, with an adjusted R-squared value of approximately 87.56%.

```
dataf <- data.frame(data)
model <- lm(price~ carat+cut+color+clarity+depth+y+z, data = data)
library(MASS)
step.model <- stepAIC(model, direction = "both", trace = FALSE)
summary(step.model)
```

```
Call:
lm(formula = price ~ cut + color + clarity + depth + y, data = data)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-2024.04  -582.36   26.43   511.83  2578.76
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -50330.66    3099.14  -16.240 < 2e-16 ***
cut           121.63      57.03    2.133  0.0342 *
color         477.39     39.46   12.100 < 2e-16 ***
clarity       777.89     46.47   16.741 < 2e-16 ***
depth        304.82     45.51    6.698 2.16e-10 ***
y            5106.73    140.37   36.381 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

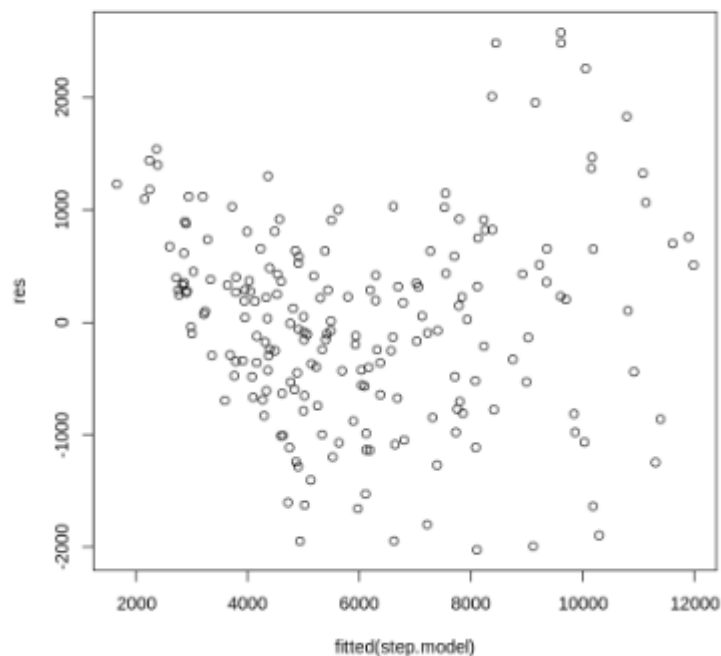
```
Residual standard error: 895.7 on 197 degrees of freedom
Multiple R-squared:  0.8786,    Adjusted R-squared:  0.8756
F-statistic: 285.2 on 5 and 197 DF,  p-value: < 2.2e-16
```

3. Create a scatter plot to examine the relationship between the fitted values and the residuals in the Linear Regression model.

The x-axis represents the predicted values, and the y-axis represents the residuals, showing the relationship between predicted and residual values. This aids in checking the completeness of the Linear Regression model, which should have residuals distributed randomly both vertically and horizontally. **If residuals scatter randomly around zero, it indicates the suitability of our Linear Regression model. However, non-random shapes or trends above or below may suggest inadequacy of the model.**

From examining the Scatter Plot, it's observed that the majority of residuals scatter above and below the zero level.

```
res <- resid(step.model)
plot(fitted(step.model), res)
```



4. Durbin-Watson Test

To check for autocorrelation in the Linear Regression model

- The lag Autocorrelation value is 0.2184589
- The D-W Statistic (Durbin-Watson Statistic) value is 1.549828
- The p-value is 0.002

The results of the Durbin-Watson test provide the D-W Statistic (Durbin-Watson Statistic), which is a value between 0 and 4, with the following interpretation:

If the D-W Statistic is closer to 2 (between 1.5 and 2.5), it indicates that there is no time-related autocorrelation in the residuals.

One of the conditions for multiple regression analysis is that the error terms must be independent. The Durbin-Watson statistic is used to check for this. If the Durbin-Watson value is close to 2, i.e., falls within the range of 1.5 to 2.5, it can be concluded that the error terms are independent.

From the data analysis, the Durbin-Watson value was found to be 1.549828, which is within the range of 1.5 to 2.5. Therefore, it can be concluded that the independent variables used in the test do not have any internal correlation.

```
install.packages("car")
library(car)
durbinWatsonTest(step.model)
```

Installing package into '/usr/local/lib/R/site-library'
(as 'lib' is unspecified)

```
lag Autocorrelation D-W Statistic p-value
1      0.2184589      1.549828  0.002
Alternative hypothesis: rho != 0
```

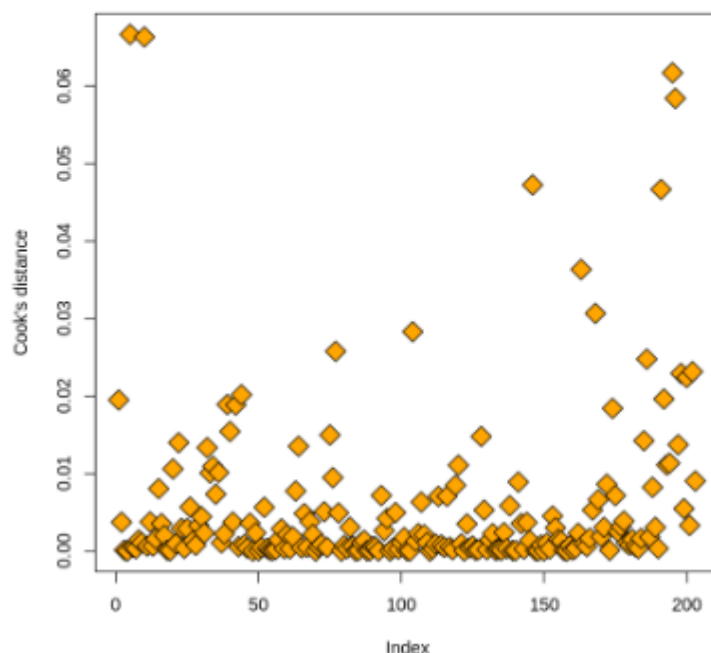
5. Cook's distance

Cook's distance is a measure used to assess how each data point in a Linear Regression model influences the estimation within the model. It is primarily used to examine data points that have a significant impact on the model or may be outliers, or data points that have the most influence on the estimation within the model. A high Cook's distance for a data point indicates that the data point may have a significant impact on the estimation within the model, requiring further consideration or detailed analysis to decide whether that data point should be included in the model or not.

When the Cook's distance falls within the range of 0.00 to 0.07, it indicates that the data point has a minimal impact on the estimation within the Linear Regression model, which is considered to have a very low impact on the estimation within the model. Overall, the model can still provide reliable results in estimating the independent variables that affect the dependent variable efficiently.

```
out <- cooks.distance(step.model)
out
1: 0.019478838679321 2: 0.00371309091246506 3: 0.00111505072207719 4: 9.80643226451709e-06 5: 0.0666335577581792 6: 0.000483196807685789 7: 0.00310287599308864 8: 0.0014267705187587 9: 0.00107236051732534 10: 0.0663018511125092 11: 0.000636554979191252 12: 0.00374975282427623 13: 0.000638841934684354 14: 0.00157199850326975 15: 0.00803637697860493 16: 0.0035351430647189 17: 0.00202045768652925 18: 0.00114226694304681 19: 5.83276267278604e-05 20: 0.0105870132843604 21: 0.000940499971335732 22: 0.013977873369013 23: 0.002922592920594 24: 0.0025810182712075 25: 0.00288372834808071 26: 0.00560503812188991 27: 0.00131534075586025 28: 0.00074650488004 29: 0.0030030215084667 30: 0.004597175010297 31: 0.0023387768803256 32: 0.01333777931637 33: 0.00057461120989 34: 0.010936531384591 35: 0.0073344423574592 36: 0.010183358646837 37: 0.00105454575832691 38: 0.0020201802021339 39: 0.01889695213048 40: 0.0154182732026542 41: 0.00373050252326186 42: 0.00058861188148635 44: 0.0201535136849239 45: 0.000808585213694767 46: 0.0003677569312629 47: 0.0035895978178943 48: 5.2210751827397e-06 49: 0.002446411484192 50: 5.8506953475142e-05 51: 0.0049680884848804 52: 0.00564370265249941 53: 0.000347699199903897 54: 0.000114181577812375 55: 8.97852135295e-05 56: 0.00012476962803511 57: 0.0014534772154045 58: 0.0002872601893929 59: 0.00031101061847488 60: 0.002101198228502 61: 0.00020844275455587 62: 0.0072381232889919 64: 0.013515246895566 65: 0.000414362588122262 66: 0.00480933190731253 67: 0.0004756767879892 68: 0.0037826239047961 69: 0.00213538898041566 70: 2.29784692370471e-06 71: 0.000548835985071108 72: 0.0008543977727338 73: 0.005066387003159 74: 0.000464941931867962 75: 0.014962115696962 76: 0.009640788541167132 77: 0.0257605466623725 78: 0.0494184004787622 79: 3.47926875943007e-07 80: 0.00072550145569183 81: 0.000328916143496973 82: 0.00306114949403815 83: 0.00066523882450138 84: 1.07175793820711e-05 85: 8.30862172491116e-05 86: 0.00107048594646434 87: 0.00137117024727578 88: 3.34307048602326e-05 89: 6.176380028185e-05 90: 0.00057991717373639 91: 0.00050821783423318 92: 2.34220070203673e-05 93: 0.007159643018947 94: 0.0026384620599366 95: 0.00419591950341283 96: 1.7729721468891e-05 97: 0.00328164052811835 98: 0.00496224336752883 99: 0.00025115383770624 100: 0.00116953121010304 101: 0.0017259353179802 102: 4.7133478031788368494e-05 103: 1.09941991918797e-05 104: 0.00887993057388667 113: 0.0070549955333666 114: 0.00077426606195787 115: 0.00048495380721768 116: 0.00706343706409355 117: 0.000490292878283787 118: 0.0001878486698131 119: 0.00847326144650538 120: 0.00110718528024455 121: 0.00092339634446022 122: 1.9329586916002e-05 123: 1.9329586916002e-05 124: 0.00349144172049764 125: 0.00039379386767723 126: 0.000436776761923412 126: 0.000104142413948941 127: 0.000174863218287601 128: 0.014768890811977 129: 0.002573214488017 130: 0.000181076249414553 131: 0.001432703039978 132: 0.0024940032378129 133: 9.575691050051e-05 134: 0.00017737839465945 135: 9.664030840720e-05 136: 0.00238080710329125 137: 0.00031369769330919 138: 0.0584327303785742 139: 7.39539432292037e-05 140: 5.6180306031089e-05 141: 0.0088991562718008 142: 0.0035903753632441 143: 0.000261874578318196 144: 0.00367789042351576 145: 0.0013628721634665 146: 0.047218352789147 147: 5.13884119602324e-05 148: 1.62563061548312e-05 149: 0.00077408968745045 150: 2.91679414749147e-05 151: 0.00099828292176137 152: 0.00027624338562502 153: 0.00451525711103029 154: 0.00299708020512059 155: 0.00150414860787534 156: 0.0012965743570287 157: 7.42191830623657e-05 158: 2.4959778811309e-06 159: 0.000677545035922 160: 0.000445507038776438 161: 0.00034488227528941 162: 0.0023806922121072 163: 0.0361516443997314 164: 0.000925904503011862 165: 0.000573893991457866 166: 0.00164670013971995 167: 0.0052792389556897 168: 0.030628149877042 169: 0.00662742894835711 170: 0.0018406533312547 171: 0.00308031693264844 172: 0.00862001771345276 173: 9.6473959682152e-05 174: 0.0184022660528 175: 0.007155631232788 176: 0.00246812903179433 177: 0.00316331139257279 178: 0.00387047230353543 179: 0.00103993926405791 180: 0.00075781619877794 181: 0.0003306031441093705 182: 0.00143690254735711 182: 0.00057433687087654 183: 0.00033103800173998 184: 0.0016219510409053 185: 0.0142254759497093 186: 0.0142254759497093 187: 0.0047462972281032 187: 0.0018649579644026 188: 0.008233195727751 189: 0.00306031441093705 190: 0.00034897470908548 191: 0.0466394610364021 192: 0.01104875045564 194: 0.011330228910217 195: 0.061675791652186 196: 0.0083894448288548 197: 0.013710712878143 198: 0.0229050468326967 199: 0.00546638928106562 200: 0.0223723646661068 201: 0.0003049649778148 202: 0.0231308346066991 203: 0.00903187244171285
```

`plot(cooks.distance(step.model), pch=23, bg='orange', cex=2, ylab="Cook's distance")`

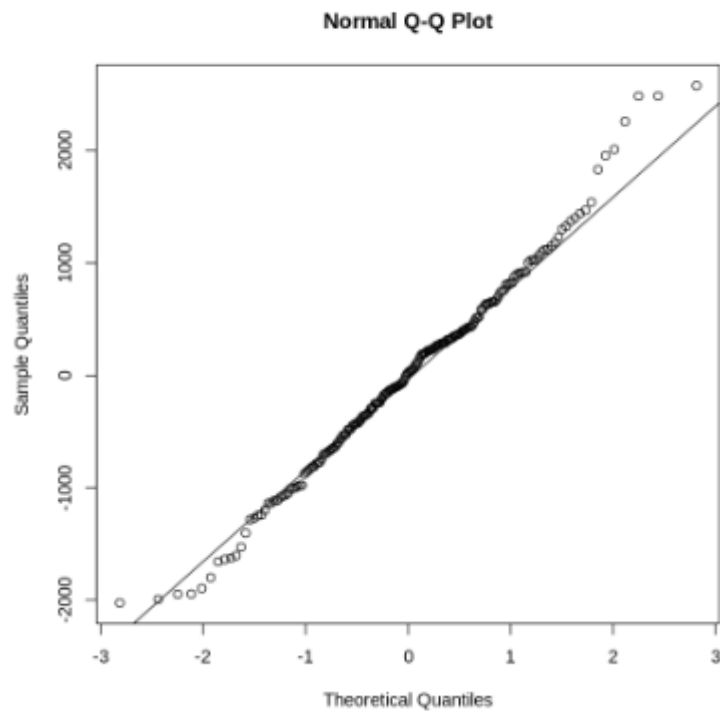


6. Quantile-Quantile (Q-Q) Plot

This is a method used to check whether the residuals of a Linear Regression model follow a normal distribution or not. It uses a Quantile-Quantile plot (QQ plot) to make predictions.

The test of residuals using a QQ plot is a method that creates a graph comparing the actual residuals with the expected values from a normal distribution. If the residuals follow a normal distribution, the data points will appear as a diagonal line with a consistent slope in the QQ plot and will be arranged in a straight line. The presence of data points on the straight line indicates a suitable distribution matching the normal distribution. However, if the residuals do not follow a normal distribution, there will be data points that deviate from the diagonal line or exhibit non-uniform clustering.

```
qqnorm(res)  
qqline(res)
```



7. Shapiro-Wilk Test

Upon evaluating the p-value, which is greater than the significance level of 0.05, we cannot reject the null hypothesis H_0 . The null hypothesis states that the data follows a normal distribution.

- W (statistic) = 0.99047
- H_0 : Data follows a normal distribution
- H_1 : Data does not follow a normal distribution
- p-value = 0.2003

Therefore, we do not have sufficient statistical evidence to reject the hypothesis that the data follows a normal distribution. Hence, we cannot conclude that the data does not follow a normal distribution.

```
shapiro.test(res)
```

Shapiro-Wilk normality test

```
data: res  
W = 0.99047, p-value = 0.2003
```

8. Testing VIF (Variance Inflation Factor)

The hypothesis used in testing VIF is:

- H_0 : There are no significant multicollinearity factors in the Linear Regression model
- H_1 : There are significant multicollinearity factors in the Linear Regression model

Therefore, when the VIF value of each independent variable is less than 10, indicating no significant multicollinearity issues in the Linear Regression model.

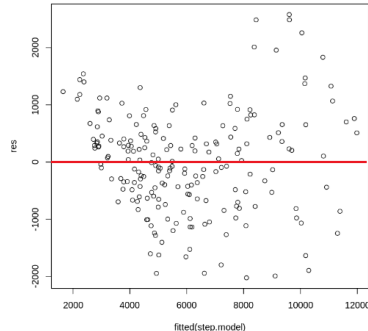
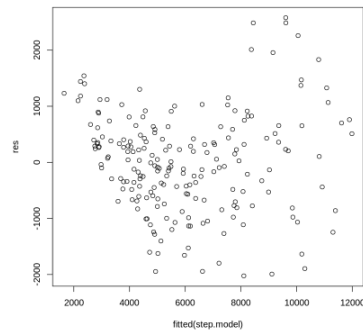
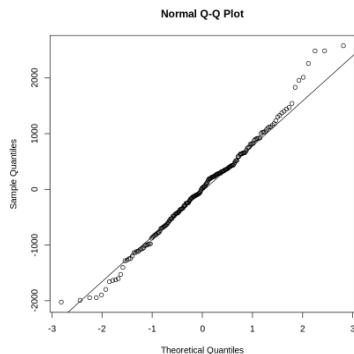
Thus, we can conclude that there are no problems regarding the interrelation between variables in the Linear Regression model, and this model can be trusted.

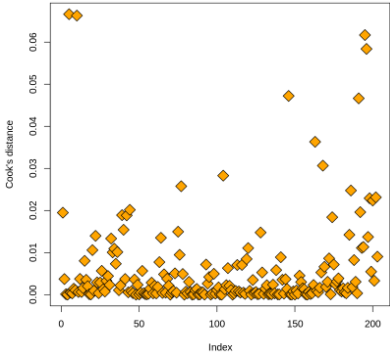
```
install.packages("car")  
library(car)  
vif(step.model)
```

Installing package into '`/usr/local/lib/R/site-library`'
(as '`lib`' is unspecified)

```
cut:      1.15260356973037 color:      1.15289761167837 clarity:      1.14244322826251 depth:      1.13794821510814 y:      1.27898539846846
```

Assumptions checking

1. Linear Model	According to the plot of residuals, they bounce randomly to ground zero. It is sufficient to suggest that the regression function is linear.									
2. Independence	The Durbin-Watson test yielded a statistic of 1.549828, if the Durbin-Watson statistic is close to 2, typically falling between 1.5 and 2.5, it suggests that the residuals are independent. In our analysis, the Durbin-Watson statistic is 1.549828, falling within the range of 1.5 to 2.5, indicating that the independent variables used in the test do not exhibit internal correlation.	<pre>install.packages("car") library(car) durbinWatsonTest(step.model)</pre> <p>Installing package into 'usr/local/lib/R/site-library' (as 'lib' is unspecified)</p> <table><tr><th>lag</th><th>Autocorrelation</th><th>D-W Statistic</th><th>p-value</th></tr><tr><td>1</td><td>0.2184589</td><td>1.549828</td><td>0.002</td></tr></table> <p>Alternative hypothesis: rho != 0</p>	lag	Autocorrelation	D-W Statistic	p-value	1	0.2184589	1.549828	0.002
lag	Autocorrelation	D-W Statistic	p-value							
1	0.2184589	1.549828	0.002							
3. Homogeneity of variance	The random pattern in the residuals vs. fitted values plot shows the random errors have constant variance.									
4. Normality	The points on the Q-Q plot fall about the straight line. That means the random errors follow a normal distribution.									

5. Outliers	There are no outliers because the Cook's distance (D_i) is less than 0.5.	 <p>A scatter plot showing Cook's distance on the y-axis (ranging from 0.00 to 0.06) against Index on the x-axis (ranging from 0 to 200). The data points are represented by orange diamonds. Most points are clustered near the bottom of the plot, with a few points reaching up to approximately 0.05. No points exceed the 0.5 threshold mentioned in the text.</p>
6. Multicollinearity	A $VIF < 5$ explains that collinearity does not affect the variable.	<pre>install.packages("car") library(car) vif(lmstep.model)</pre> <p>Installing package into '/usr/local/lib/R/site-library' (as 'lib' is unspecified)</p> <p>call: <code>1.35260256873037 color: 1.15286761167637 clarity: 1.14244522826251 depth: 1.13754621518614 p: 1.27898538646846</code></p>

Summary and Analysis

The report on analyzing and predicting diamond prices using regression emphasizes the significance of regression in forecasting data behavior across various domains. This analysis focuses on creating a statistical model to predict diamond prices using diverse characteristics such as weight, cut quality, color, clarity, and size. The research aims to develop an appropriate statistical model to enhance understanding and accurate prediction of diamond prices, highlighting the importance of regression in solving mathematical problems and aiding decision-making in the diamond industry.

BLUE (Best Linear Unbiased Estimators) Analysis:

- Unbiasedness: The model meets the assumption of unbiasedness as the regression coefficients are estimated without bias.
- Linearity: The residuals scatter randomly around zero, indicating that the regression function is linear.
- Efficiency: Given that the model satisfies the assumptions, the estimators are efficient and achieve the smallest variance among all unbiased estimators.

ANOVA & t-test Assumptions:

- Independence: The Durbin-Watson statistic falling within the range of 1.5 to 2.5 indicates that the residuals are independent, meeting the assumption for ANOVA and t-tests.
- Homogeneity of Variance: The variance of the residual value compared to the filled plot value shows that random error has constant variance.
- Normality: Normal distribution of residual values can be observed in the Q-Q plot, indicating that the random error values follow a normal distribution.

Model Effectiveness:

- Outliers: The absence of outliers, indicated by Cook's distance being less than 0.5, ensures that influential data points are not skewing the results.
- Multicollinearity: If the VIF value is less than 5, it indicates no issue with multicollinearity, suggesting that the variables do not have high interrelationships.

The results of the analysis and prediction of diamond prices using regression indicate that the 5 independent variables (cut, color, clarity, depth, y(width)) significantly influence diamond prices. The linear regression model developed has a strong ability to explain the data, with an adjusted R-squared value of approximately 87.56%. Additionally, the p-value of the model is less than the standard statistical significance level of 0.05, allowing the rejection of the null hypothesis (H0: independent variables have no effect on diamond prices) and acceptance of the alternative hypothesis (H1: at least one independent variable affects diamond prices), suggesting a response of diamond prices to the input independent variables.