

[PREDICTING QUALITY OF WINE]

R ABISHEK

BSC DATA SCIENCE AND ANALYTICS

Introduction

Multiple linear regression is a statistical method used to analyze the relationship between two or more independent variables and a single dependent variable. It is an extension of simple linear regression, which considers only one independent variable. In multiple linear regression, we aim to model the linear relationship between the independent variables (X_1, X_2, \dots, X_p) and the dependent variable (Y) using the equation:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

where:

- (Y) is the dependent variable (also known as the response or outcome variable).
- (X_1, X_2, \dots, X_p) are the independent variables (also known as predictors or features).
- $(\beta_0, \beta_1, \beta_2, \dots, \beta_p)$ are the regression coefficients, representing the slope of the relationship between each independent variable and the dependent variable.
- (ϵ) is the error term, representing the difference between the observed and predicted values of the dependent variable.

The goal of multiple linear regression is to estimate the values of the regression coefficients (β) 's that minimize the sum of squared errors (SSE) between the observed and predicted values of the dependent variable. This is typically done using the method of least squares.

Aim

The aim of this analysis is to develop a predictive model using multiple linear regression to estimate the quality rating of wine based on its chemical properties. By leveraging a dataset containing information on red and white variants of the Portuguese "Vinho Verde" wine, including various chemical attributes and quality ratings, we seek to understand the relationship between these factors and wine quality.

Through the application of multiple linear regression, we aim to:

1. Identify the key chemical properties that significantly influence the quality rating of wine.
2. Build a predictive model that accurately estimates wine quality based on its chemical composition.
3. Evaluate the performance of the predictive model using appropriate metrics and validation techniques.
4. Provide insights and recommendations based on the findings to stakeholders in the wine industry, such as producers, distributors, and consumers.

By achieving these objectives, we can contribute to a better understanding of the factors contributing to wine quality and provide practical tools for decision-making and quality control within the wine industry. Moreover, the predictive model developed in this analysis

can serve as a valuable asset for optimizing production processes, enhancing product quality, and meeting consumer preferences.

Dataset Overview

The dataset under consideration contains information on red and white variants of the Portuguese "Vinho Verde" wine. It comprises various chemical properties of the wines, such as acidity levels, pH, alcohol content, and more. Additionally, each wine sample is associated with a quality rating, providing a measure of its perceived quality.

The dataset is structured in tabular form, with rows representing individual wine samples and columns representing different attributes or features. Typically, the first few rows of the dataset provide an insight into its structure and contents, including the names of the variables and the range of values present.

Understanding the dataset's overview is essential before conducting any analysis. It helps in identifying the scope of the data, determining the variables of interest, and planning the analysis approach accordingly. Additionally, knowing the dataset's source and any specific data collection processes can provide context for interpreting the results of subsequent analyses.

Dataset Source

The dataset originates from the UCI Machine Learning Repository, a widely-used repository for machine learning datasets. The specific dataset we are working with contains information on red and white variants of the Portuguese "Vinho Verde" wine.

The UCI Machine Learning Repository serves as a valuable resource for researchers, practitioners, and students in the field of machine learning and data science. It hosts a diverse collection of datasets spanning various domains, including healthcare, finance, education, and more.

Datasets in the UCI repository are often curated, cleaned, and made publicly available for academic and research purposes. They may originate from research studies, industry sources, or public data sources. Detailed documentation and information about each dataset are typically provided to facilitate their usage in research and experimentation.

By leveraging datasets from reputable sources like the UCI Machine Learning Repository, researchers and practitioners can access high-quality data for analysis, model development, and evaluation. Moreover, the availability of such datasets promotes transparency, reproducibility, and collaboration in the field of machine learning and data science.

Importing CSV file

Input:

```
# Import the merged dataset
winequality <- read.csv("/Users/m2air/Downloads/WineData/winequality.csv", header = TRUE)
```

Dataset Dimension

Input:

```
# Check the dimensions of the dataset
dim(winequality)
```

Output:

```
> dim(winequality)
[1] 6497  13
```

Result:

No. of Columns – 13

No. of Rows – 6497

Dataset Variables

Input:

```
# Display column names
names(winequality)
```

Output:

```
> names(winequality)
[1] "fixed.acidity"      "volatile.acidity"   "citric.acid"
[4] "residual.sugar"     "chlorides"          "free.sulfur.dioxide"
[7] "total.sulfur.dioxide" "density"            "pH"
[10] "sulphates"          "alcohol"            "quality"
[13] "type"
```

Result:

1. Fixed acidity: The total concentration of acids that do not evaporate readily.
2. Volatile acidity: The concentration of volatile acids that contribute to the wine's sourness or vinegar-like taste.
3. Citric acid: The concentration of citric acid in the wine, which can impart freshness and flavour.
4. Residual sugar: The residual sugar content in the wine after fermentation, contributing to sweetness.
5. Chlorides: The concentration of chlorides in the wine, affecting taste and stability.

6. Free sulphur dioxide: The concentration of free sulphur dioxide, acting as an antioxidant and antimicrobial agent.
7. Total sulphur dioxide: The total concentration of sulphur dioxide, including both free and bound forms.
8. Density: The density of the wine, influenced by alcohol and sugar content.
9. pH: The pH level of the wine, measuring acidity or alkalinity.
10. Sulphates: The concentration of sulphates, contributing to aroma and flavour stability.
11. Alcohol: The alcohol content of the wine, affecting body, flavour, and quality.
12. Quality: The quality rating assigned to each wine sample.
13. Type: Indicates the type of wine (red or white).

These variables collectively provide information about the chemical properties and quality characteristics of the wine samples, with the "type" variable indicating whether the wine is red or white.

Exploratory Data Analysis

Structure of the Dataset:

`str(winequality)`

Output:

```
> str(winequality)
'data.frame':  6497 obs. of  13 variables:
 $ fixed.acidity      : num  7.4 7.8 7.8 11.2 7.4 7.4 7.9 7.3 7.8 7.5 ...
 $ volatile.acidity   : num  0.7 0.88 0.76 0.28 0.7 0.66 0.6 0.65 0.58 0.5 ...
 $ citric.acid        : num  0 0 0.04 0.56 0 0 0.06 0 0.02 0.36 ...
 $ residual.sugar     : num  1.9 2.6 2.3 1.9 1.9 1.8 1.6 1.2 2 6.1 ...
 $ chlorides          : num  0.076 0.098 0.092 0.075 0.076 0.075 0.069 0.065 0.073 0.071
 ...
 $ free.sulfur.dioxide : num  11 25 15 17 11 13 15 15 9 17 ...
 $ total.sulfur.dioxide: num  34 67 54 60 34 40 59 21 18 102 ...
 $ density             : num  0.998 0.997 0.997 0.998 0.998 ...
 $ pH                  : num  3.51 3.2 3.26 3.16 3.51 3.51 3.3 3.39 3.36 3.35 ...
 $ sulphates           : num  0.56 0.68 0.65 0.58 0.56 0.56 0.46 0.47 0.57 0.8 ...
 $ alcohol             : num  9.4 9.8 9.8 9.8 9.4 9.4 9.4 10 9.5 10.5 ...
 $ quality             : int  5 5 5 6 5 5 5 7 7 5 ...
 $ type                : chr  "red" "red" "red" "red" ...
```

Summary Statistics of the Dataset:

`summary(winequality)`

Output:

`> summary(winequality)`

fixed.acidity	volatile.acidity	citric.acid	residual.sugar
Min. : 3.800	Min. : 0.0800	Min. : 0.0000	Min. : 0.600
1st Qu.: 6.400	1st Qu.: 0.2300	1st Qu.: 0.2500	1st Qu.: 1.800
Median : 7.000	Median : 0.2900	Median : 0.3100	Median : 3.000
Mean : 7.215	Mean : 0.3397	Mean : 0.3186	Mean : 5.443
3rd Qu.: 7.700	3rd Qu.: 0.4000	3rd Qu.: 0.3900	3rd Qu.: 8.100
Max. : 15.900	Max. : 1.5800	Max. : 1.6600	Max. : 65.800

chlorides	free.sulfur.dioxide	total.sulfur.dioxide	density
Min. : 0.00900	Min. : 1.00	Min. : 6.0	Min. : 0.9871
1st Qu.: 0.03800	1st Qu.: 17.00	1st Qu.: 77.0	1st Qu.: 0.9923
Median : 0.04700	Median : 29.00	Median : 118.0	Median : 0.9949
Mean : 0.05603	Mean : 30.53	Mean : 115.7	Mean : 0.9947
3rd Qu.: 0.06500	3rd Qu.: 41.00	3rd Qu.: 156.0	3rd Qu.: 0.9970
Max. : 0.61100	Max. : 289.00	Max. : 440.0	Max. : 1.0390

pH	sulphates	alcohol	quality	type
Min. : 2.720	Min. : 0.2200	Min. : 8.00	Min. : 3.000	Length:6497
1st Qu.: 3.110	1st Qu.: 0.4300	1st Qu.: 9.50	1st Qu.: 5.000	Class :character
Median : 3.210	Median : 0.5100	Median : 10.30	Median : 6.000	Mode :character
Mean : 3.219	Mean : 0.5313	Mean : 10.49	Mean : 5.818	
3rd Qu.: 3.320	3rd Qu.: 0.6000	3rd Qu.: 11.30	3rd Qu.: 6.000	
Max. : 4.010	Max. : 2.0000	Max. : 14.90	Max. : 9.000	

First few Columns:

`head(winequality)`

Output:

`> head(winequality)`

	fixed.acidity	volatile.acidity	citric.acid	residual.sugar	chlorides
1	7.4	0.70	0.00	1.9	0.076
2	7.8	0.88	0.00	2.6	0.098
3	7.8	0.76	0.04	2.3	0.092
4	11.2	0.28	0.56	1.9	0.075
5	7.4	0.70	0.00	1.9	0.076
6	7.4	0.66	0.00	1.8	0.075

	free.sulfur.dioxide	total.sulfur.dioxide	density	pH	sulphates	alcohol	quality
1	11	34	0.9978	3.51	0.56	9.4	5
2	25	67	0.9968	3.20	0.68	9.8	5
3	15	54	0.9970	3.26	0.65	9.8	5
4	17	60	0.9980	3.16	0.58	9.8	6
5	11	34	0.9978	3.51	0.56	9.4	5
6	13	40	0.9978	3.51	0.56	9.4	5

	type
1	red
2	red
3	red
4	red
5	red
6	red

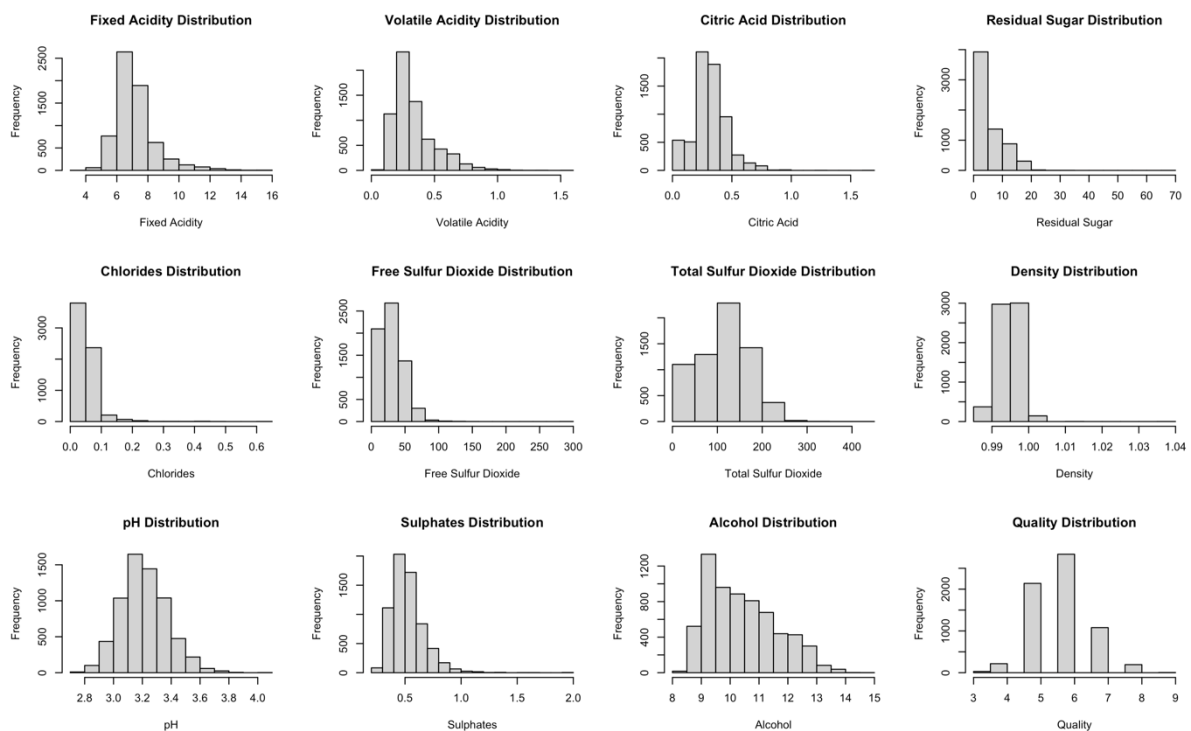
Variable Analysis

- Analyse each variable to understand its distribution.

Input:

```
# Display histograms for numerical variables (excluding 'type')
par(mfrow=c(3, 4)) # Set up a 3x4 grid for plots
hist(winequality$fixed.acidity, main = "Fixed Acidity Distribution", xlab = "Fixed Acidity")
hist(winequality$volatile.acidity, main = "Volatile Acidity Distribution", xlab = "Volatile Acidity")
hist(winequality$citric.acid, main = "Citric Acid Distribution", xlab = "Citric Acid")
hist(winequality$residual.sugar, main = "Residual Sugar Distribution", xlab = "Residual Sugar")
hist(winequality$chlorides, main = "Chlorides Distribution", xlab = "Chlorides")
hist(winequality$free.sulfur.dioxide, main = "Free Sulfur Dioxide Distribution", xlab = "Free Sulfur Dioxide")
hist(winequality$total.sulfur.dioxide, main = "Total Sulfur Dioxide Distribution", xlab = "Total Sulfur Dioxide")
hist(winequality$density, main = "Density Distribution", xlab = "Density")
hist(winequality$pH, main = "pH Distribution", xlab = "pH")
hist(winequality$sulphates, main = "Sulphates Distribution", xlab = "Sulphates")
hist(winequality$alcohol, main = "Alcohol Distribution", xlab = "Alcohol")
hist(winequality$quality, main = "Quality Distribution", xlab = "Quality")
```

Output:



Result:

- The fixed acidity of most wines falls between 7 and 12.
- The volatile acidity of most wines falls between 0.3 and 0.7.
- The citric acid of most wines falls between 0 and 0.3.
- The residual sugar of most wines falls between 2 and 8.
- The free sulphur dioxide of most wines falls between 50 and 150.

Outlier Analysis

- Analyse each variable to identify Outliers.

Input:

```
# Set up the layout for plots
```

```
par(mfrow = c(3, 4))
```

```
# Create boxplots for numerical variables by wine type
```

```
boxplot(fixed.acidity ~ type, data = winequality, main = "Fixed Acidity by Wine Type", xlab = "Wine Type", ylab = "Fixed Acidity")
```

```
boxplot(volatile.acidity ~ type, data = winequality, main = "Volatile Acidity by Wine Type", xlab = "Wine Type", ylab = "Volatile Acidity")
```

```
boxplot(citric.acid ~ type, data = winequality, main = "Citric Acid by Wine Type", xlab = "Wine Type", ylab = "Citric Acid")
```

```
boxplot(residual.sugar ~ type, data = winequality, main = "Residual Sugar by Wine Type", xlab = "Wine Type", ylab = "Residual Sugar")
```

```
boxplot(chlorides ~ type, data = winequality, main = "Chlorides by Wine Type", xlab = "Wine Type", ylab = "Chlorides")
```

```
boxplot(free.sulfur.dioxide ~ type, data = winequality, main = "Free Sulfur Dioxide by Wine Type", xlab = "Wine Type", ylab = "Free Sulfur Dioxide")
```

```
boxplot(total.sulfur.dioxide ~ type, data = winequality, main = "Total Sulfur Dioxide by Wine Type", xlab = "Wine Type", ylab = "Total Sulfur Dioxide")
```

```
boxplot(density ~ type, data = winequality, main = "Density by Wine Type", xlab = "Wine Type", ylab = "Density")
```

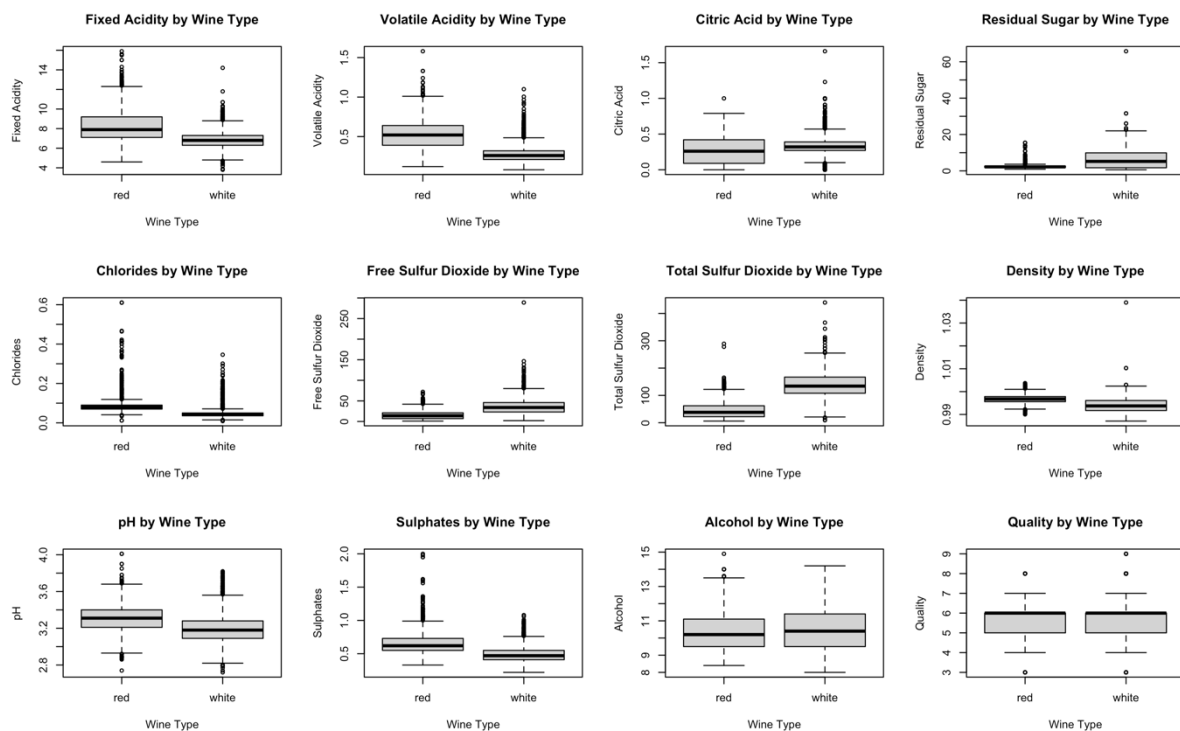
```
boxplot(pH ~ type, data = winequality, main = "pH by Wine Type", xlab = "Wine Type", ylab = "pH")
```

```
boxplot(sulphates ~ type, data = winequality, main = "Sulphates by Wine Type", xlab = "Wine Type", ylab = "Sulphates")
```

```
boxplot(alcch2 ~ type, data = winequality, main = "Alcohol by Wine Type", xlab = "Wine Type", ylab = "Alcohol")
```

```
boxplot(quality ~ type, data = winequality, main = "Quality by Wine Type", xlab = "Wine Type", ylab = "Quality")
```


Output:



Result:

Some variables has outliers which may affect the accuracy of the model.

Addressing Outliers

To address outliers in all variables shortly using R, we can use a simple approach such as Winsorization or trimming. Here's how you can Winsorize the dataset:

```
# Function to Winsorize a vector
winsorize <- function(x, p = 0.05) {
  q <- quantile(x, probs = c(p, 1 - p), na.rm = TRUE)
  x[x < q[1]] <- q[1]
  x[x > q[2]] <- q[2]
  return(x)
}

# Apply Winsorization to all numerical variables except 'type'
numerical_vars <- winequality[, sapply(winequality, is.numeric)]
winequality[, names(numerical_vars)] <- sapply(numerical_vars, winsorize)

# Display summary statistics after Winsorization
summary(winequality)
```

This code defines a Winsorization function ``winsorize`` and applies it to all numerical variables in the dataset except for the "type" variable. Winsorization replaces extreme values with the values at the specified quantiles (5th and 95th percentiles by default). After Winsorization, you can check the summary statistics of the dataset to see if the outliers have been mitigated. Adjust the parameters of the ``winsorize`` function as needed.

Correlation Matrix

A correlation matrix shows how strongly pairs of variables are related to each other. It gives a quick overview of the relationships between variables, with values ranging from -1 to 1, where:

- 1 means a perfect positive correlation (as one variable increases, the other also increases),
- -1 means a perfect negative correlation (as one variable increases, the other decreases), and
- 0 means no correlation between the variables.

The correlation matrix helps identify which variables are positively or negatively associated with each other, aiding in understanding the data's patterns and relationships.

Input:

```
# Extract the relevant quantitative predictor variables
predictor_variables <- winequality[, c("fixed.acidity", "volatile.acidity", "citric.acid",
"residual.sugar", "chlorides", "free.sulfur.dioxide", "total.sulfur.dioxide", "density", "pH",
"sulphates", "alcohol")]

# Calculate the correlation matrix
correlation_matrix <- cor(predictor_variables)

# Print the correlation matrix
print("Correlation Matrix:")
print(correlation_matrix)
```

Output:

	fixed.acidity	volatile.acidity	citric.acid	residual.sugar	chlorides
fixed.acidity	1.0000000	0.24400409	0.300628668	-0.1180773	0.42359572
volatile.acidity	0.2440041	1.00000000	-0.399890818	-0.2150553	0.53446875
citric.acid	0.3006287	-0.39989082	1.000000000	0.1411473	-0.08157548
residual.sugar	-0.1180773	-0.21505531	0.141147319	1.0000000	-0.17485998
chlorides	0.4235957	0.53446875	-0.081575483	-0.1748600	1.00000000
free.sulfur.dioxide	-0.2980244	-0.38331798	0.149732525	0.4496950	-0.30470255
total.sulfur.dioxide	-0.3290517	-0.43777244	0.199190101	0.5154710	-0.40830630
density	0.4482767	0.28414497	0.076216432	0.5294762	0.52969428
pH	-0.2338213	0.25900697	-0.332723698	-0.2813684	0.16241437
sulphates	0.2895584	0.25449774	0.039680065	-0.1993764	0.37932423
alcohol	-0.1003201	-0.04936782	0.007364332	-0.3638218	-0.32699823
	free.sulfur.dioxide	total.sulfur.dioxide	density	pH	
fixed.acidity	-0.2980244	-0.32905165	0.44827672	-0.23382132	
volatile.acidity	-0.3833180	-0.43777244	0.28414497	0.25900697	
citric.acid	0.1497325	0.19919010	0.07621643	-0.33272370	
residual.sugar	0.4496950	0.51547100	0.52947622	-0.28136844	
chlorides	-0.3047026	-0.40830630	0.52969428	0.16241437	
free.sulfur.dioxide	1.0000000	0.73953561	0.02940080	-0.16408097	
total.sulfur.dioxide	0.7395356	1.00000000	0.02015523	-0.25244294	
density	0.0294008	0.02015523	1.00000000	0.01890145	
pH	-0.1640810	-0.25244294	0.01890145	1.00000000	
sulphates	-0.2118424	-0.30090070	0.26717583	0.24293124	
alcohol	-0.1820038	-0.26363189	-0.70490597	0.10881586	
	sulphates	alcohol			
fixed.acidity	0.289558435	-0.100320093			
volatile.acidity	0.254497742	-0.049367818			
citric.acid	0.039680065	0.007364332			
residual.sugar	-0.199376419	-0.363821751			
chlorides	0.379324233	-0.326998228			
free.sulfur.dioxide	-0.211842365	-0.182003843			
total.sulfur.dioxide	-0.300900698	-0.263631888			
density	0.267175827	-0.704905970			
pH	0.242931239	0.108815863			

Checking Multicollinearity

We check for multicollinearity to ensure that the predictor variables in our regression model are not highly correlated with each other. High multicollinearity can lead to issues such as:

1. Unreliable Estimates: It can make estimates of the regression coefficients unstable and less precise.
2. Difficulty in Interpretation: Highly correlated predictors make it challenging to interpret the effect of each predictor variable on the outcome variable independently.
3. Difficulty in Identifying Important Variables: Multicollinearity can mask the true relationships between predictors and the outcome variable, making it difficult to identify which predictors are truly important.

By checking for multicollinearity, we ensure that our regression model produces reliable and interpretable results.

Input:

```
# Check for multicollinearity issue
max_correlation <- max(abs(correlation_matrix))
if (max_correlation < 0.9) {
  print("There is no multicollinearity issue in this model.")
} else {
  print("Multicollinearity may be a problem in this model.")
}
```

Output:

```
[1] "Multicollinearity may be a problem in this model."
```

Variance Inflation Factors

VIF measures how much the variance of a regression coefficient is inflated due to multicollinearity. It tells us how much a predictor variable is affected by the correlation with other predictors. High VIF values indicate strong multicollinearity, which can make regression coefficients unreliable.

Input:

```
# Compute Variance Inflation Factors (VIF) manually
vif_values <- sapply(1:ncol(predictor_variables), function(i) {
  lm_out <- lm(predictor_variables[, i] ~ ., data = predictor_variables[, -i])
  1 / (1 - summary(lm_out)$r.squared)
})

# Print Variance Inflation Factors (VIF)
print("Variance Inflation Factors (VIF):")
print(vif_values)
```

Output:

```
> print(vif_values)
 [1]  3.897291  2.053305  1.612539  7.045335  2.680453  2.302134  3.144361 15.898995  2.246263
[10]  1.520742  4.684512
```

Rechecking Multicollinearity using High Correlation

Input:

```
# Calculate the correlation matrix
correlation_matrix <- cor(predictor_variables)

# Find highly correlated variable pairs
correlated_pairs <- which(correlation_matrix > 0.9 & correlation_matrix < 1, arr.ind = TRUE)

# Print highly correlated variable pairs
print("Highly Correlated Variable Pairs:")
print(correlated_pairs)

# Remove one variable from each highly correlated pair
if (length(correlated_pairs) > 0) {
  vars_to_remove <- unique(c(correlated_pairs[, 1], correlated_pairs[, 2]))
  predictor_variables <- predictor_variables[, -vars_to_remove]
  print("Variables removed due to multicollinearity:")
  print(colnames(predictor_variables))
} else {
  print("There are no highly correlated variables.")
}
```

Output:

```
[1] "There are no highly correlated variables."
```

Multiple Linear Regression

Multiple linear regression, in shorter terms, is a statistical method used to understand the relationship between one dependent variable and two or more independent variables. It helps to predict the value of the dependent variable based on the values of the independent variables.

In essence, it's like fitting a straight line (or a plane in higher dimensions) to the data to find the best-fitting model that explains the relationship between the variables.

The model equation looks like this:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Where:

- Y is the dependent variable,
- (X_1, X_2, \dots, X_n) are the independent variables,

- $(\beta_0, \beta_1, \beta_2, \dots, \beta_n)$ are the coefficients (slope) representing the effect of each independent variable on the dependent variable,
- ϵ is the error term representing the difference between the observed and predicted values, and
- The goal is to estimate the coefficients $(\beta_0, \beta_1, \beta_2, \dots, \beta_n)$ that best fit the data.

In simpler terms, multiple linear regression helps us understand how changes in one or more independent variables are associated with changes in the dependent variable. It allows us to make predictions based on the relationships observed in the data.

Model Building:

```
# Fit the multiple linear regression model
model <- lm(quality ~ ., data = data)
```

Model Evaluation:

Input:

```
# Summary of the model
summary(model)
```

Output:

Call:

```
lm(formula = quality ~ ., data = data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-3.7796	-0.4671	-0.0444	0.4561	3.0211

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.048e+02	1.414e+01	7.411	1.42e-13	***
fixed.acidity	8.507e-02	1.576e-02	5.396	7.05e-08	***
volatile.acidity	-1.492e+00	8.135e-02	-18.345	< 2e-16	***
citric.acid	-6.262e-02	7.972e-02	-0.786	0.4322	
residual.sugar	6.244e-02	5.934e-03	10.522	< 2e-16	***
chlorides	-7.573e-01	3.344e-01	-2.264	0.0236	*
free.sulfur.dioxide	4.937e-03	7.662e-04	6.443	1.25e-10	***
total.sulfur.dioxide	-1.403e-03	3.237e-04	-4.333	1.49e-05	***
density	-1.039e+02	1.434e+01	-7.248	4.71e-13	***
pH	4.988e-01	9.058e-02	5.506	3.81e-08	***
sulphates	7.217e-01	7.624e-02	9.466	< 2e-16	***
alcohol	2.227e-01	1.807e-02	12.320	< 2e-16	***
typewhite	-3.613e-01	5.675e-02	-6.367	2.06e-10	***

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

Residual standard error: 0.7331 on 6484 degrees of freedom

Multiple R-squared: 0.2965, Adjusted R-squared: 0.2952

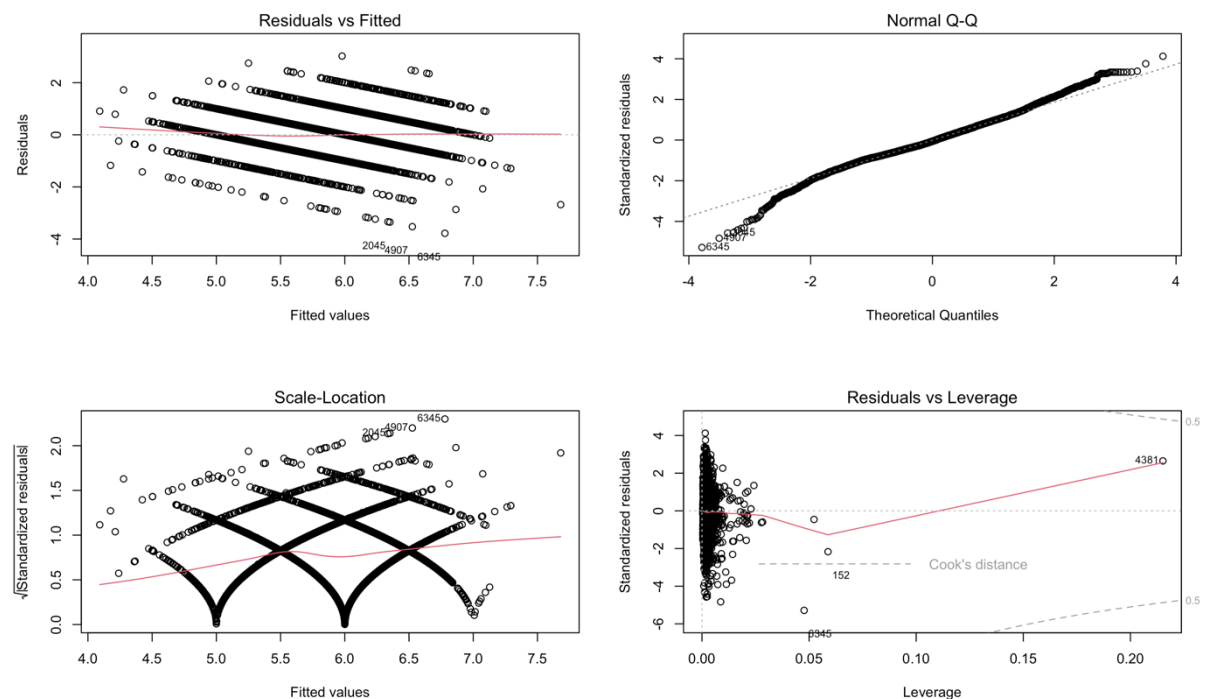
F-statistic: 227.8 on 12 and 6484 DF, p-value: < 2.2e-16

Plots:

Input:

```
# Assess model assumptions  
par(mfrow = c(2, 2)) # Arrange plots in a 2x2 grid  
plot(model)
```

Output:



Conclusions:

Residuals vs Fitted

- The x-axis shows the fitted values, which are the predicted values for the response variable based on the model.
- The y-axis shows the residuals, which are the difference between the actual values of the response variable and the fitted values.
- In a perfect model, the residuals would be randomly scattered around zero on the y-axis. This would indicate that there is no relationship between the fitted values and the residuals.
- The residuals appear to be scattered around zero, but there is a slight trend. This suggests that there may be a weak relationship between the fitted values and the residuals. However, it is difficult to say for sure from this graph alone.

Normal Q-Q

- ☐ A normal Q-Q plot is a graphical technique for comparing two probability distributions.
- ☐ In this case, the normal Q-Q plot is comparing the distribution of the residuals to a normal distribution.
- ☐ If the residuals are normally distributed, the points on the normal Q-Q plot will fall close to a straight line.
- ☐ The points on the normal Q-Q plot deviate from a straight line, particularly at the tails. This suggests that the residuals are not normally distributed.

Scale-Location

- ☐ I couldn't find any information about what a "Scale-Location" plot is used for.
- ☐ The x-axis appears to show the fitted values.
- ☐ The y-axis appears to show the standardized residuals.
- ☐ Standardized residuals are a way of measuring the residuals in units of standard deviations.
- ☐ In a normal distribution, about 68% of the data points will fall within 1 standard deviation of the mean, and 95% of the data points will fall within 2 standard deviations of the mean.
- ☐ It appears that most of the standardized residuals fall within 2 standard deviations of the mean. However, there are a few outliers that fall outside of this range.

Residuals vs Leverage

- ☐ Leverage is a measure of how much influence a particular data point has on the fitted model.
- ☐ Data points with high leverage are more likely to be outliers.
- ☐ The x-axis of the graph you sent shows the leverage of each data point.
- ☐ The y-axis shows the residuals.
- ☐ In a perfect model, there would be no relationship between the leverage and the residuals.
- ☐ There appears to be a slight trend, with the residuals increasing as the leverage increases. This suggests that there may be some outliers that are having a large influence on the fitted model.

Prediction:

```
# Make predictions
```

```
predictions <- predict(model, data)
```

Prediction Accuracy:

Input:

```
# Assess prediction accuracy
```

```
accuracy <- sqrt(mean((data$quality - predictions)^2))
```

```
print(paste("Root Mean Squared Error (RMSE):", round(accuracy, 2)))
```

Output:

```
[1] "Root Mean Squared Error (RMSE): 0.73"
```

Extracting Coefficients

Input:

```
coefficients <- coef(model)
```

```
# Print coefficients
```

```
print(coefficients)
```

Output:

(Intercept)	fixed.acidity	volatile.acidity	citric.acid
1.047518e+02	8.506605e-02	-1.492406e+00	-6.262113e-02
residual.sugar	chlorides	free.sulfur.dioxide	total.sulfur.dioxide
6.243733e-02	-7.572508e-01	4.937045e-03	-1.402677e-03
density	pH	sulphates	alcohol
-1.039096e+02	4.987558e-01	7.217444e-01	2.226697e-01
typewhite			
-3.613317e-01			

Regression Equation

Based on the coefficients , we can construct the regression equation for predicting wine quality based on the variables in the WineQuality dataset:

$$\begin{aligned} \text{Quality} = & 104.75 + 0.0851 \times \text{Fixed Acidity} - 1.4924 \times \text{Volatile Acidity} - \\ & 0.0626 \times \text{Citric Acid} + 0.0624 \times \text{Residual Sugar} - 0.7573 \times \text{Chlorides} + 0.0049 \times \\ & \text{Free Sulfur Dioxide} - 0.0014 \times \text{Total Sulfur Dioxide} - 103.91 \times \text{Density} + \\ & 0.4988 \times \text{pH} + 0.7217 \times \text{Sulphates} + 0.2227 \times \text{Alcohol} - 0.3613 \times \text{Type White} \end{aligned}$$

This equation suggests how each predictor variable contributes to the predicted wine quality. For example, an increase in fixed acidity is associated with a slight increase in wine quality, while an increase in volatile acidity is associated with a decrease in wine quality.

Dataset Reference Link

<https://archive.ics.uci.edu/dataset/186/wine+quality>