

NetID: qiyangw3

Qiyang Wang

2024-12-01

Contents

Abstract	1
Model Choice: Full Model and Null Model with Stepwise Selection	3
transformation and diagnostic plots	6
outlier and influential points	12
colinearity	16
LASSO, Ridge, and Elastic Net Regression	17
Conclusion	21

Abstract

This study explores the relationship between calorie expenditure and various factors such as age, gender, workout type, session duration, and heart rate among gym enthusiasts. Using regression analysis and random forest modeling, we developed predictive models to estimate calorie burn based on these variables. The analysis highlights key factors influencing calorie expenditure and provides insights into optimizing workout routines for individual fitness goals. Our findings offer practical guidance for fitness enthusiasts and professionals to design personalized and efficient exercise programs. ## analysis the dataset

```
## 'data.frame':   973 obs. of  15 variables:
##  $ Age                : int  56 46 32 25 38 56 36 40 28 28 ...
##  $ Gender              : chr  "Male" "Female" "Female" "Male" ...
##  $ Weight..kg.         : num  88.3 74.9 68.1 53.2 46.1 ...
##  $ Height..m.          : num  1.71 1.53 1.66 1.7 1.79 1.68 1.72 1.51 1.94 1.84 ...
##  $ Max_BPM             : int  180 179 167 190 188 168 174 189 185 169 ...
##  $ Avg_BPM             : int  157 151 122 164 158 156 169 141 127 136 ...
##  $ Resting_BPM         : int  60 66 54 56 68 74 73 64 52 64 ...
##  $ Session_Duration..hours. : num  1.69 1.3 1.11 0.59 0.64 1.59 1.49 1.27 1.03 1.08 ...
##  $ Calories_Burned      : num  1313 883 677 532 556 ...
##  $ Workout_Type        : chr  "Yoga" "HIIT" "Cardio" "Strength" ...
##  $ Fat_Percentage       : num  12.6 33.9 33.4 28.8 29.2 15.5 21.3 30.6 28.9 29.7 ...
##  $ Water_Intake..liters. : num  3.5 2.1 2.3 2.1 2.8 2.7 2.3 1.9 2.6 2.7 ...
##  $ Workout_Frequency..days.week.: int  4 4 4 3 3 5 3 3 4 3 ...
##  $ Experience_Level     : int  3 2 2 1 1 3 2 2 2 1 ...
##  $ BMI                 : num  30.2 32 24.7 18.4 14.4 ...

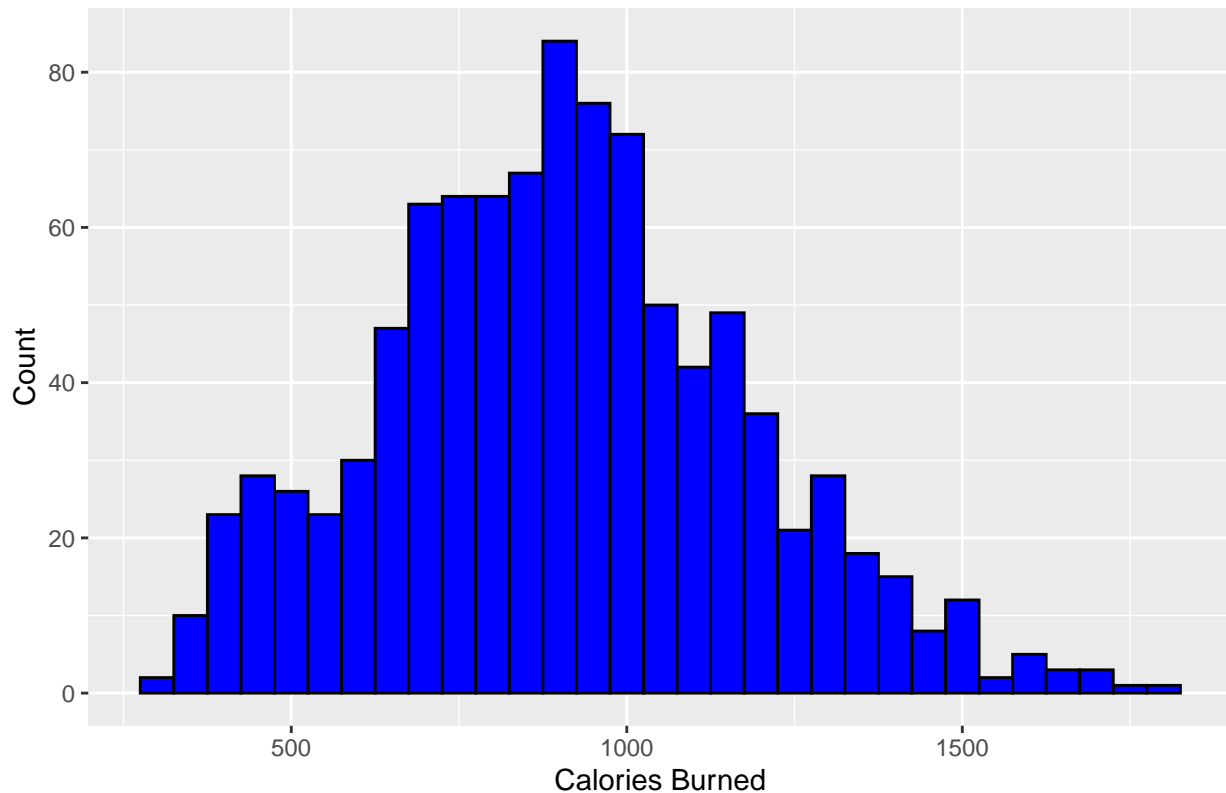
##      Age      Gender      Weight..kg.      Height..m.
##  Min.   :18.00   Length:973   Min.    : 40.00   Min.    :1.500
##  1st Qu.:28.00   Class :character   1st Qu.: 58.10   1st Qu.:1.620
##  Median :40.00   Mode  :character   Median : 70.00   Median :1.710
##  Mean   :38.68                Mean   : 73.85   Mean   :1.723
```

```

## 3rd Qu.:49.00          3rd Qu.: 86.00    3rd Qu.:1.800
## Max.   :59.00          Max.    :129.90   Max.    :2.000
##      Max_BPM          Avg_BPM          Resting_BPM      Session_Duration..hours.
## Min.   :160.0        Min.    :120.0    Min.    :50.00   Min.    :0.500
## 1st Qu.:170.0        1st Qu.:131.0    1st Qu.:56.00   1st Qu.:1.040
## Median :180.0        Median :143.0    Median :62.00   Median :1.260
## Mean   :179.9        Mean    :143.8    Mean    :62.22   Mean    :1.256
## 3rd Qu.:190.0        3rd Qu.:156.0    3rd Qu.:68.00   3rd Qu.:1.460
## Max.   :199.0        Max.    :169.0    Max.    :74.00   Max.    :2.000
## Calories_Burned      Workout_Type          Fat_Percentage   Water_Intake..liters.
## Min.    : 303.0      Length:973          Min.    :10.00   Min.    :1.500
## 1st Qu.: 720.0      Class :character    1st Qu.:21.30   1st Qu.:2.200
## Median : 893.0      Mode  :character    Median :26.20   Median :2.600
## Mean    : 905.4                      Mean    :24.98   Mean    :2.627
## 3rd Qu.:1076.0                      3rd Qu.:29.30   3rd Qu.:3.100
## Max.    :1783.0                      Max.    :35.00   Max.    :3.700
## Workout_Frequency..days.week. Experience_Level      BMI
## Min.    :2.000                      Min.    :1.00   Min.    :12.32
## 1st Qu.:3.000                      1st Qu.:1.00   1st Qu.:20.11
## Median :3.000                      Median :2.00   Median :24.16
## Mean    :3.322                      Mean    :1.81   Mean    :24.91
## 3rd Qu.:4.000                      3rd Qu.:2.00   3rd Qu.:28.56
## Max.    :5.000                      Max.    :3.00   Max.    :49.84

```

Distribution of Calories Burned



```

##                               Age                               Weight..kg.
##                               -0.154678760                     0.095443473
##                               Height..m.                       Max_BPM

```

```
##                0.086348051                0.002090016
##                Avg_BPM                Resting_BPM
##                0.339658667                0.016517951
##      Session_Duration..hours.                Calories_Burned
##                0.908140376                1.000000000
##                Fat_Percentage                Water_Intake..liters.
##                -0.597615248                0.356930683
## Workout_Frequency..days.week.                Experience_Level
##                0.576150125                0.694129448
##                BMI
##                0.059760826
```

After loading the dataset and performing some initial exploration, we observed the following:

The histogram indicates that `Calories_Burned` has a unimodal distribution with a slight skew to the right. Most values fall within the range of 800 to 1400. This suggests that most observations are clustered in this range, but there are a few higher values that might require further attention for outliers or special cases.

By calculating the correlation between `Calories_Burned` and other numeric variables, we identified the strength of their linear relationships. Strongly correlated variables (positive or negative) might have more predictive power for our target variable, while weak correlations might indicate limited predictive value.

Next Step: fit a full model to understand the relationships between `Calories_Burned` and other variables in the dataset, we will start with a full model that includes all predictors

Model Choice: Full Model and Null Model with Stepwise Selection

To develop a predictive model for `Calories_Burned`, both a **full model** and a **null model** were used as starting points for **stepwise selection** to find the optimal set of predictors.

Steps:

1. **Data Split:** The dataset was divided into 70% training data and 30% testing data for model training and evaluation.
2. **Full Model:**
 - The full model included all predictors to assess their collective contribution to explaining `Calories_Burned`.
 - This model served as the upper limit for the stepwise selection process.
3. **Null Model:**
 - The null model only included the intercept, assuming no predictors were significant.
 - This model served as the lower limit for stepwise selection.
4. **Stepwise Selection:**
 - **BIC:** Stepwise selection based on the Bayesian Information Criterion identified a parsimonious model with fewer predictors by strongly penalizing model complexity.
 - **AIC:** Stepwise selection based on the Akaike Information Criterion provided a more flexible approach, potentially retaining more predictors for better predictive performance.

Conclusion: Using both the full and null models in stepwise selection ensures a systematic approach to identifying the optimal subset of predictors. BIC favors simpler models, while AIC allows for slightly more complexity, providing a balance between interpretability and predictive accuracy.

```
##      (Intercept) Session_Duration..hours.                Avg_BPM
##      -820.835603                715.587845                6.369233
##      GenderMale                Age
```

88.949297 -3.495094

Model Formula-stepwise selection based on BIC

The final linear regression model for predicting `Calories_Burned` is as follows:

$$\text{Calories_Burned} = -820.8356 + 715.5878 \cdot \text{Session_Duration_hours} + 6.3692 \cdot \text{Avg_BPM} + 88.9493 \cdot \text{GenderMale} - 3.4951 \cdot \text{Age}$$

Model Coefficients The table below shows the coefficients estimated by the regression:

Predictor	Coefficient	Std. Error	t-value	p-value	Significance
Intercept	-820.8356	16.8076	-48.84	< 2e-16	***
Session_Duration..hours.	715.5878	4.4116	162.21	< 2e-16	***
Avg_BPM	6.3692	0.1046	60.89	< 2e-16	***
GenderMale	88.9493	3.0310	29.35	< 2e-16	***
Age	-3.4951	0.1243	-28.13	< 2e-16	***

Key Metrics

- **Residual Standard Error (RSE):** 39.42
- **Multiple R-squared:** 0.9795
- **Adjusted R-squared:** 0.9794
- **F-statistic:** 8085 on 4 and 676 degrees of freedom
- **p-value:** < 2.2e-16

Key Metrics

- **R-squared:** 0.9795
- **Adjusted R-squared:** 0.9794
- **Residual Standard Error:** 39.42
- **F-statistic (p-value):** 8085 (< 2.2e-16)

Interpretation

1. **Session Duration:** Each additional hour burns **715.6 more calories**.
2. **Avg_BPM:** A 1-unit increase burns **6.37 additional calories**.
3. **Gender (Male):** Males burn **88.95 more calories** than females.
4. **Age:** Each additional year decreases calorie burn by **3.5 calories**.

Conclusion This model explains 97.95% of the variance in `Calories_Burned` and identifies session duration, average BPM, gender, and age as significant predictors. It provides a reliable tool for estimating calorie expenditure and optimizing workout strategies.

Model Formula-stepwise selection based on AIC

##	(Intercept)	Session_Duration..hours.	Avg_BPM
##	-881.9136772	715.9314060	6.3561892
##	GenderMale	Age	Height..m.
##	85.1221175	-3.4842386	25.1815303
##	Resting_BPM		
##	0.3322244		

Model Formula The AIC-optimized linear regression model for predicting **Calories_Burned** is:

$$\begin{aligned}\text{Calories_Burned} = & -881.9137 + 715.9314 \cdot \text{Session_Duration_hours} + \\ & 6.3562 \cdot \text{Avg_BPM} + 85.1221 \cdot \text{GenderMale} - \\ & 3.4842 \cdot \text{Age} + 25.1815 \cdot \text{Height_m} + \\ & 0.3322 \cdot \text{Resting_BPM}\end{aligned}$$

Coefficients

Predictor	Coefficient	Std. Error	p-value	Significance
Intercept	-881.9137	32.2916	< 2e-16	***
Session_Duration..hours.	715.9314	4.4036	< 2e-16	***
Avg_BPM	6.3562	0.1048	< 2e-16	***
GenderMale	85.1221	3.7317	< 2e-16	***
Age	-3.4842	0.1241	< 2e-16	***
Height..m.	25.1815	14.4701	0.0823	.
Resting_BPM	0.3322	0.2121	0.1178	

Key Results

- **R-squared:** 0.9797
- **Residual Standard Error:** 39.32
- **F-statistic:** 5417 ($p < 2.2\text{e-}16$)

Significant Predictors

1. **Session Duration:** Each additional hour increases calories burned by **715.93**.
2. **Avg BPM:** Each unit increase adds **6.36 calories**.
3. **Gender (Male):** Males burn **85.12 more calories** than females.
4. **Age:** Each additional year reduces calories burned by **3.48**.

Conclusion The AIC-optimized model explains **97.97%** of the variance in calories burned. Key predictors like session duration, heart rate, gender, and age provide actionable insights for optimizing fitness plans.

Determine the best model And then we can use the BIC and AIC value to determine the model

```
## The AIC of stepwise_model_bic: 6943.891
## The BIC of stepwise_model_bic: 6971.033
## The AIC of stepwise_model_aic: 6942.609
## The BIC of stepwise_model_aic: 6978.797
```

Results:

- **Stepwise Model (BIC):**
 - AIC: 6943.891
 - BIC: 6971.033
- **Stepwise Model (AIC):**
 - AIC: 6942.609
 - BIC: 6978.797

Why Choose AIC? We selected the AIC-based model because it prioritizes **predictive accuracy**. While the BIC-based model slightly simplifies the model by penalizing complexity more heavily, the AIC-based model is more suited for maximizing prediction performance. The lower AIC value (6942.609) supports the choice of the AIC model.

Theoretical Background: AIC and BIC

1. **AIC (Akaike Information Criterion)** The AIC evaluates a model's goodness of fit while penalizing complexity. The formula is:

$$\text{AIC} = -2 \cdot \text{Log-Likelihood} + 2k$$

Where: - Log-Likelihood: Measures how well the model fits the data. - k : Number of estimated parameters in the model.

- **Interpretation:** Lower AIC values indicate a better trade-off between fit and complexity. AIC is more flexible and allows slightly more parameters to improve predictive accuracy.
2. **BIC (Bayesian Information Criterion)** The BIC similarly balances fit and complexity but applies a stronger penalty for additional predictors. The formula is:

$$\text{BIC} = -2 \cdot \text{Log-Likelihood} + k \cdot \log(n)$$

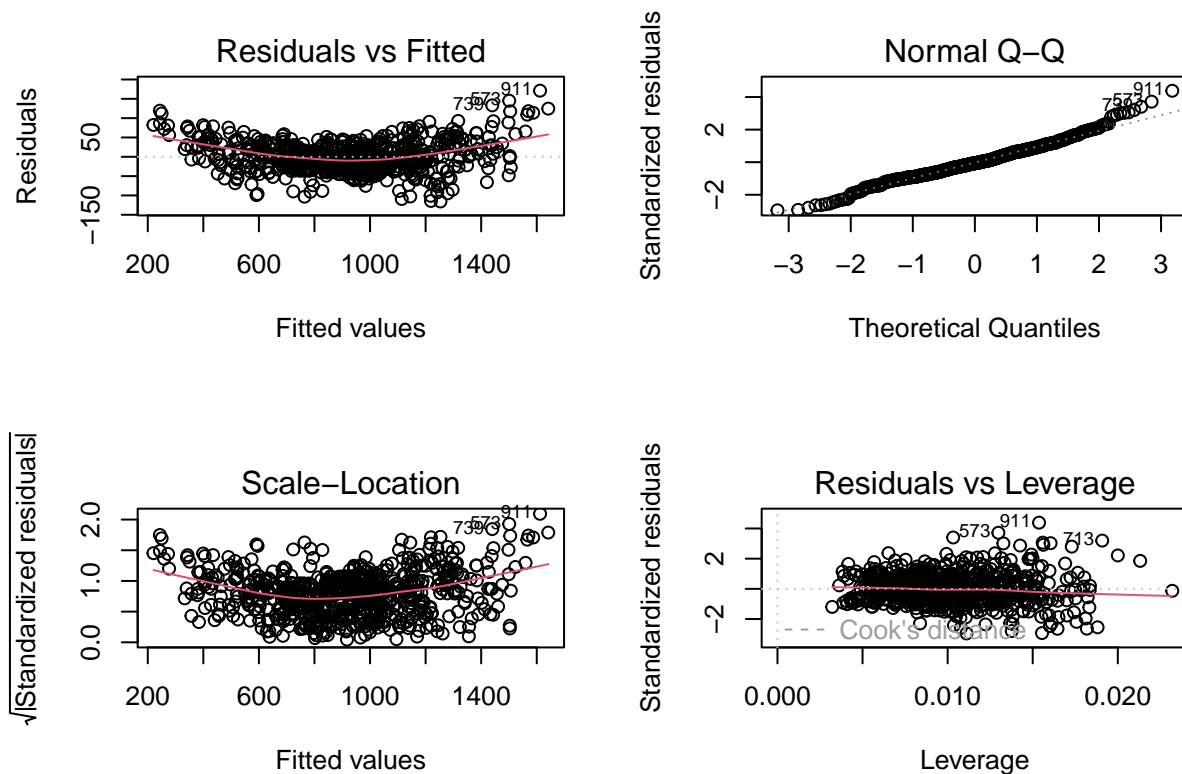
Where: - n : Number of observations. - k : Number of estimated parameters in the model.

- **Interpretation:** Lower BIC values indicate a better model. BIC tends to favor simpler models and is more conservative than AIC.

transformation and diagnostic plots

we want a better prediction, so we decide to use the AIC model.

##	(Intercept)	Session_Duration..hours.	Avg_BPM
##	-881.9136772	715.9314060	6.3561892
##	GenderMale	Age	Height..m.
##	85.1221175	-3.4842386	25.1815303
##	Resting_BPM		
##	0.3322244		



The diagnostic plots help us assess the assumptions of the linear regression model. We can check for linearity, homoscedasticity, normality of residuals, and influential points. If any of these assumptions are violated, we may need to address them before interpreting the model results.

Non-Linearity: The “Residuals vs Fitted” plot shows a curved pattern, suggesting a non-linear relationship between the predictors and the response.

Non-Normality: The “Normal Q-Q” plot shows deviations at the tails, indicating the residuals are not perfectly normal.

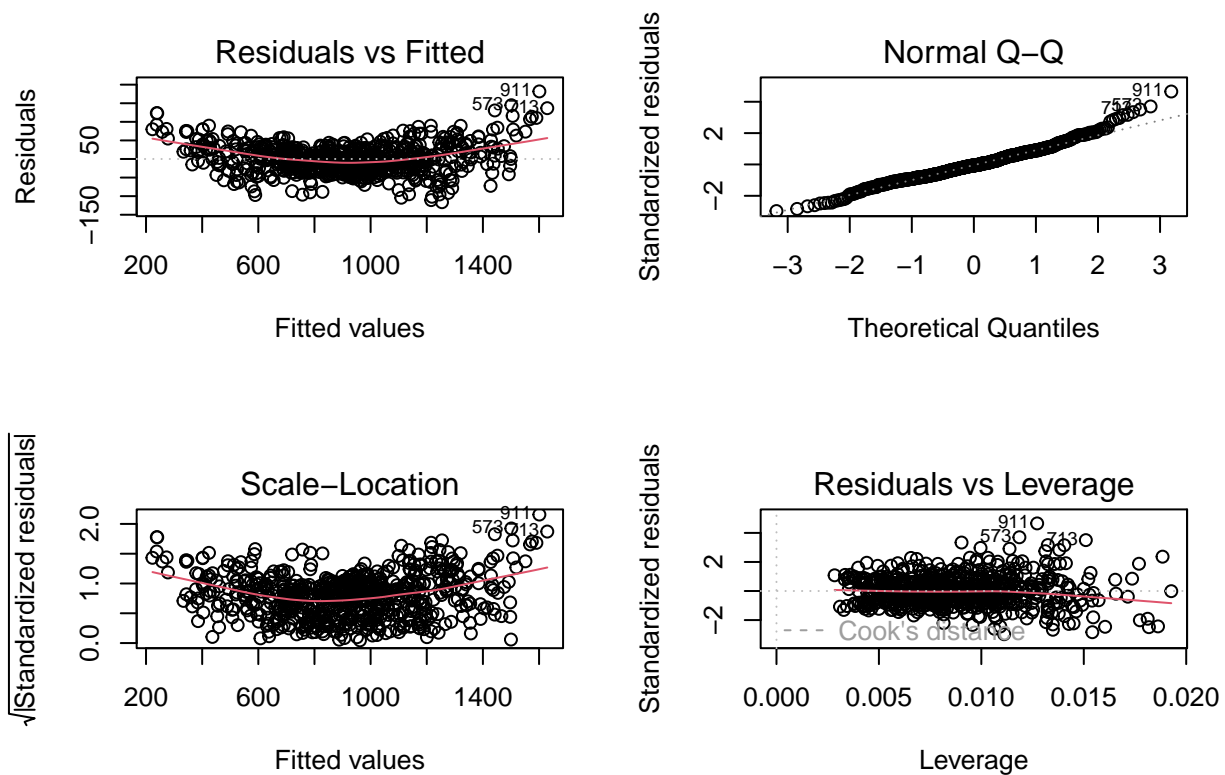
Heteroscedasticity: The “Scale-Location” plot shows a slight pattern, indicating that the variance of residuals may not be constant.

To address these issues, we can try transforming the predictors or the response variable to linearize relationships, stabilize variance, and normalize residuals.

Add log transformation to the model

```
##
## Call:
## lm(formula = Calories_Burned ~ Session_Duration..hours. + log(Avg_BPM) +
##     Gender + Age + log(Resting_BPM), data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -116.248  -26.843   -2.657   23.422  181.829
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)          -4508.068      87.959 -51.252  <2e-16 ***
## Session_Duration..hours.  715.524      4.402 162.532  <2e-16 ***
## log(Avg_BPM)           911.802     15.034  60.649  <2e-16 ***
## GenderMale             88.797       3.024  29.366  <2e-16 ***
## Age                   -3.509       0.124 -28.306  <2e-16 ***
## log(Resting_BPM)       18.963     13.011   1.457   0.145
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 39.32 on 675 degrees of freedom
## Multiple R-squared:  0.9797, Adjusted R-squared:  0.9795
## F-statistic: 6500 on 5 and 675 DF, p-value: < 2.2e-16
```



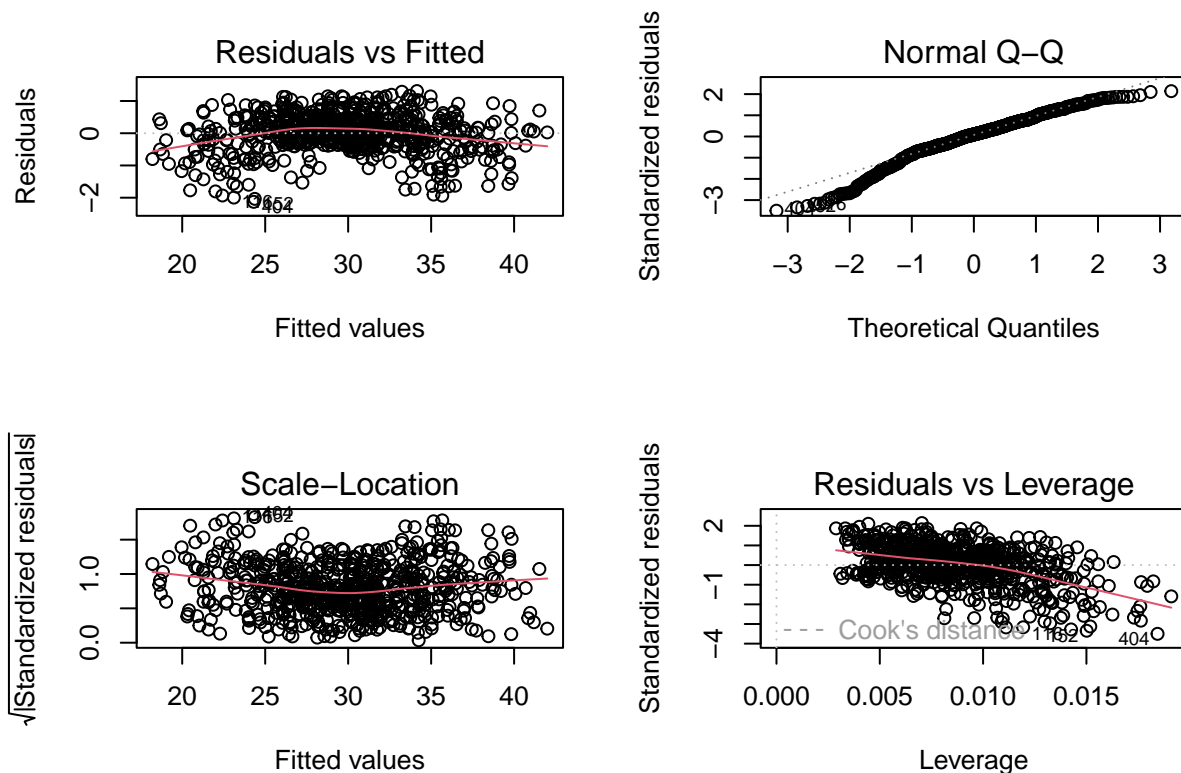
The log transformation improved model validity by addressing key issues with linearity, normality, and variance consistency. However, further investigation of influential points and potential additional transformations is recommended to fully optimize the model.

Add square root transformation to the model

```
##
## Call:
## lm(formula = sqrt(Calories_Burned) ~ Session_Duration..hours. +
##     Avg_BPM + Gender + Age + Resting_BPM, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.12821 -0.32294  0.06018  0.41394  1.30784
```



```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    7.633e-01  3.211e-01   2.377  0.0177 *
## Session_Duration..hours. 1.213e+01  6.870e-02 176.613 <2e-16 ***
## Avg_BPM         1.053e-01  1.635e-03  64.417 <2e-16 ***
## GenderMale      1.446e+00  4.719e-02  30.654 <2e-16 ***
## Age            -5.649e-02  1.935e-03 -29.199 <2e-16 ***
## Resting_BPM     9.879e-05  3.307e-03   0.030  0.9762
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6137 on 675 degrees of freedom
## Multiple R-squared:  0.9825, Adjusted R-squared:  0.9824
## F-statistic: 7591 on 5 and 675 DF, p-value: < 2.2e-16
```



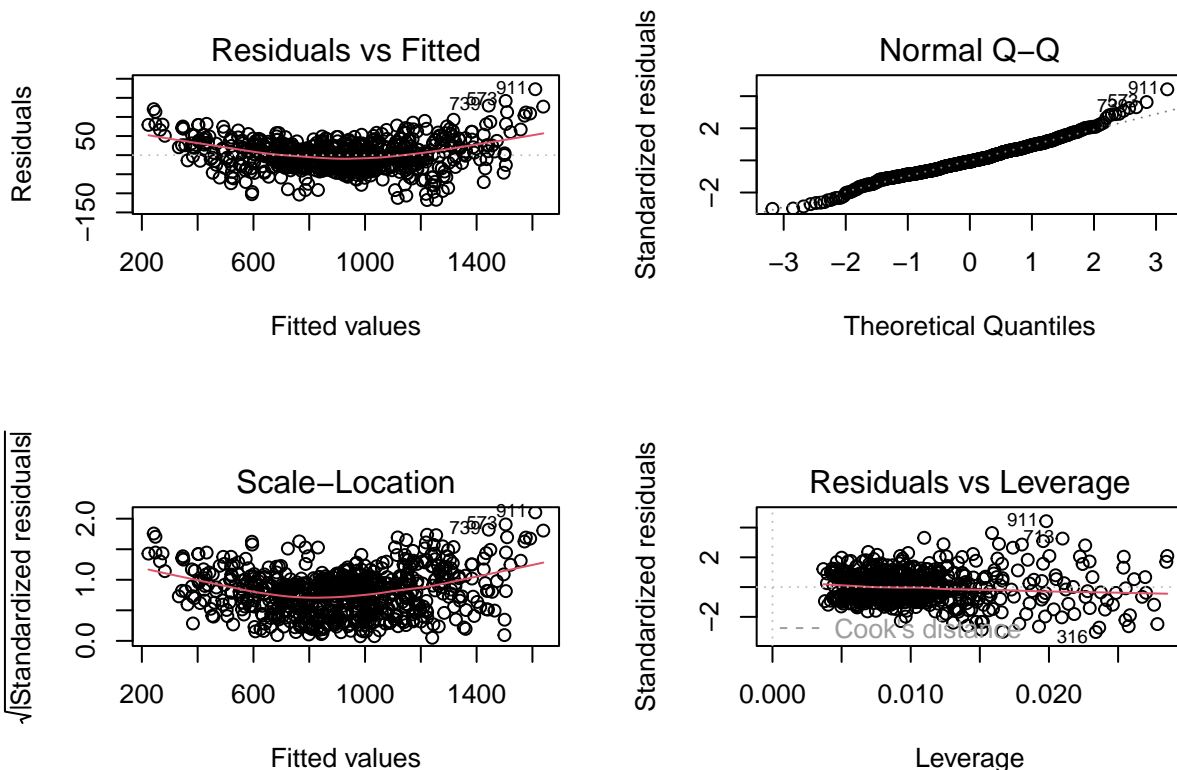
The square root transformation showed some improvement in addressing: - **Linearity:** Curvature is reduced but not eliminated. - **Normality:** Residuals are closer to normality. - **Homoscedasticity:** Variance of residuals is more stable.

However, the improvements are limited compared to the log-transformed model. To capture the remaining non-linear relationships, adding **quadratic terms** or exploring other transformations may be necessary. Additionally, influential points should be reviewed for potential removal or robust modeling.

Add quadratic terms to the model

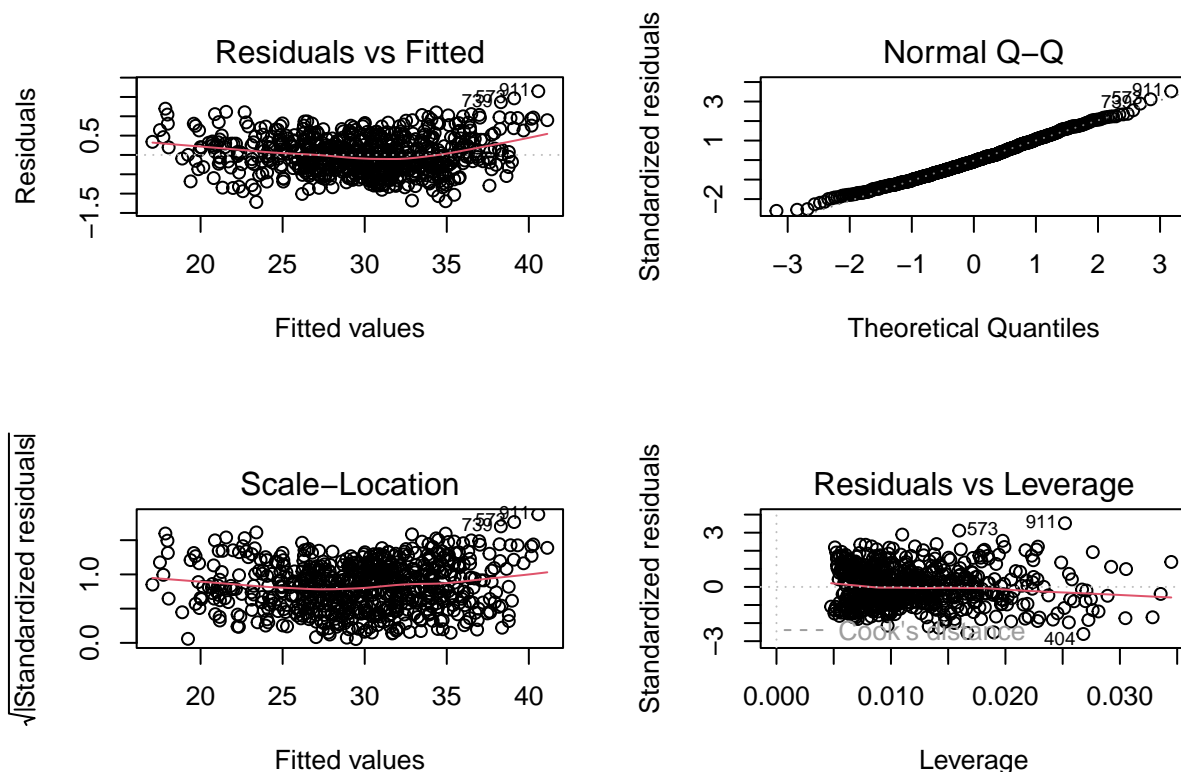
```
##
## Call:
```

```
## lm(formula = Calories_Burned ~ Session_Duration..hours. + I(Session_Duration..hours.^2) +
##     Avg_BPM + Gender + Age + Resting_BPM, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -117.495  -25.879   -2.684    25.199   172.628
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -830.5902    25.3769  -32.730  <2e-16 ***
## Session_Duration..hours.    701.9309    25.8902   27.112  <2e-16 ***
## I(Session_Duration..hours.^2)    5.4142     9.9998    0.541    0.588
## Avg_BPM         6.3581     0.1052   60.457  <2e-16 ***
## GenderMale     88.9679     3.0308   29.354  <2e-16 ***
## Age          -3.4956     0.1242  -28.141  <2e-16 ***
## Resting_BPM     0.3104     0.2126    1.460    0.145
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 39.4 on 674 degrees of freedom
## Multiple R-squared:  0.9796, Adjusted R-squared:  0.9794
## F-statistic: 5394 on 6 and 674 DF, p-value: < 2.2e-16
```



```
##
## Call:
## lm(formula = sqrt(Calories_Burned) ~ poly(Session_Duration..hours.,
```

```
##      2) + poly(Avg_BPM, 2) + Gender + Age + Resting_BPM, data = train_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.21683 -0.33839 -0.02982  0.31541  1.64839
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      31.146584    0.171168  181.965 < 2e-16 ***
## poly(Session_Duration..hours., 2)1  108.577001    0.474859  228.651 < 2e-16 ***
## poly(Session_Duration..hours., 2)2 -10.049548    0.475136  -21.151 < 2e-16 ***
## poly(Avg_BPM, 2)1      39.147791    0.476745   82.115 < 2e-16 ***
## poly(Avg_BPM, 2)2     -1.692661    0.474319   -3.569 0.000384 ***
## GenderMale          1.425752    0.036450   39.116 < 2e-16 ***
## Age                -0.056822    0.001495  -38.014 < 2e-16 ***
## Resting_BPM         0.002505    0.002556    0.980 0.327487
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4738 on 673 degrees of freedom
## Multiple R-squared:  0.9896, Adjusted R-squared:  0.9895
## F-statistic: 9161 on 7 and 673 DF, p-value: < 2.2e-16
```



The addition of quadratic terms significantly improved model diagnostics: - **poly_model1** addressed key issues like curvature, normality, and heteroscedasticity. - **poly_model2**, with both square root transformation and polynomial terms, achieved slightly better normality and variance stability.

Given the improved fit, `poly_model2` may be preferred for prediction due to its ability to handle non-linear relationships and improve residual diagnostics. Further refinement could include addressing influential points or testing interaction terms for additional predictors.

outlier and influential points

Outliers and influential points can have a significant impact on the model's performance. We can identify these points using diagnostic plots and leverage techniques like Cook's distance to detect influential observations.

Influential Points and High Leverage Points Analysis

To evaluate the impact of specific data points on the regression model, we analyzed **Cook's Distance** and **Leverage Values**. These diagnostics help identify points that may unduly influence the model's estimates or predictions.

Leverage Values Definition: Leverage measures how far an observation's predictor values are from the mean of the predictors. High leverage points can disproportionately affect the fit of the model.

Formula:

$$h_{ii} = \mathbf{x}_i (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i^\top$$

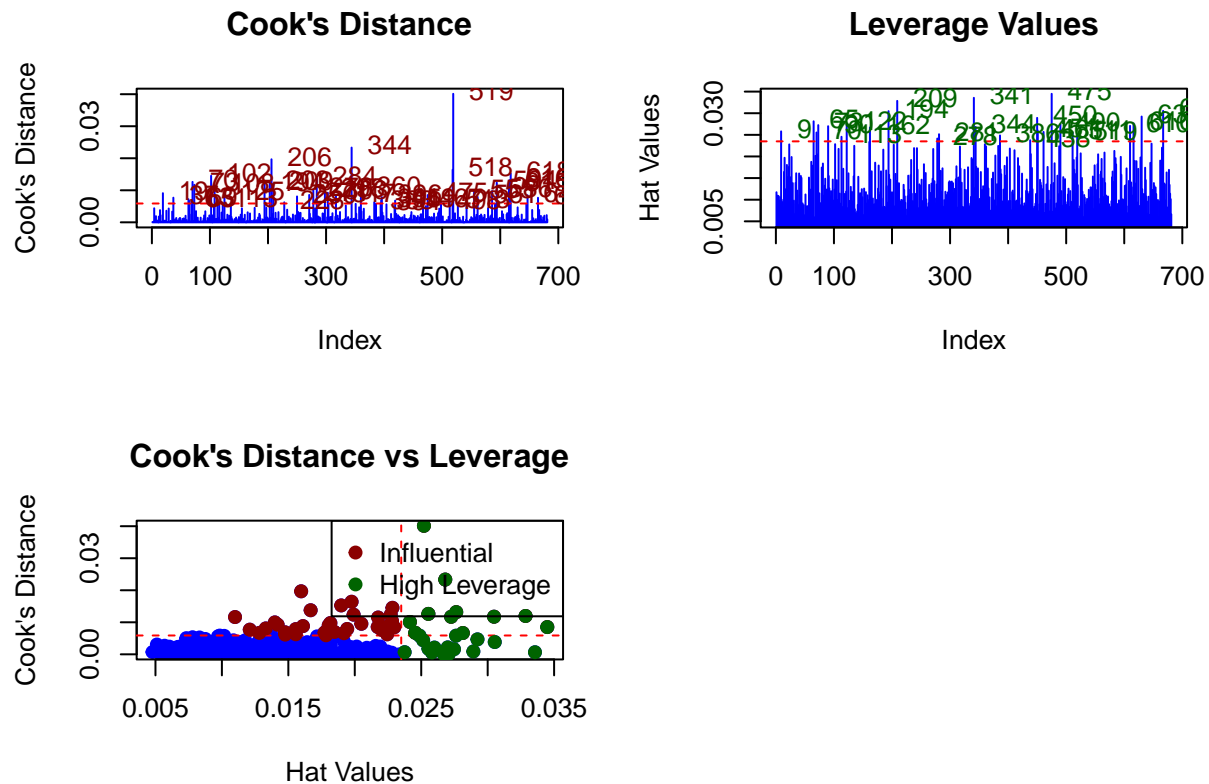
Where: - \mathbf{x}_i : The row vector of predictor values for observation i . - \mathbf{X} : The matrix of all predictor values.

Threshold: Points with $h_{ii} > \frac{2p}{n}$ are considered high leverage.

Table 3: Summary of Influential and High Leverage Points

Type	Index
Influential Points	19
Influential Points	37
Influential Points	63
Influential Points	65
Influential Points	69
Influential Points	70
Influential Points	73
Influential Points	102
Influential Points	108
Influential Points	113
Influential Points	125
Influential Points	194
Influential Points	203
Influential Points	206
Influential Points	209
Influential Points	228
Influential Points	250
Influential Points	271
Influential Points	278
Influential Points	284
Influential Points	298
Influential Points	303
Influential Points	317
Influential Points	344
Influential Points	360

Type	Index
Influential Points	383
Influential Points	394
Influential Points	395
Influential Points	396
Influential Points	419
Influential Points	467
Influential Points	475
Influential Points	498
Influential Points	513
Influential Points	518
Influential Points	519
Influential Points	555
Influential Points	561
Influential Points	586
Influential Points	594
Influential Points	615
Influential Points	616
Influential Points	618
Influential Points	647
Influential Points	665
High Leverage Points	9
High Leverage Points	65
High Leverage Points	70
High Leverage Points	73
High Leverage Points	90
High Leverage Points	113
High Leverage Points	122
High Leverage Points	162
High Leverage Points	194
High Leverage Points	209
High Leverage Points	278
High Leverage Points	281
High Leverage Points	341
High Leverage Points	344
High Leverage Points	386
High Leverage Points	438
High Leverage Points	450
High Leverage Points	451
High Leverage Points	461
High Leverage Points	475
High Leverage Points	490
High Leverage Points	511
High Leverage Points	519
High Leverage Points	610
High Leverage Points	616
High Leverage Points	630
High Leverage Points	667



Key Findings:

1. Influential Points (Cook's Distance):

- Observations such as **519, 518, 344, 278, 284** exceed the Cook's Distance threshold, indicating significant influence on the model.

2. High Leverage Points:

- Points such as **9, 70, 113, 209, 519** have high leverage, meaning they are far from the average predictor values.

3. Overlap:

- Observations like **519, 344, and 278** are both influential and high leverage, making them critical for further review.

Implications:

- These points may distort the model, leading to biased coefficients or poor predictions.

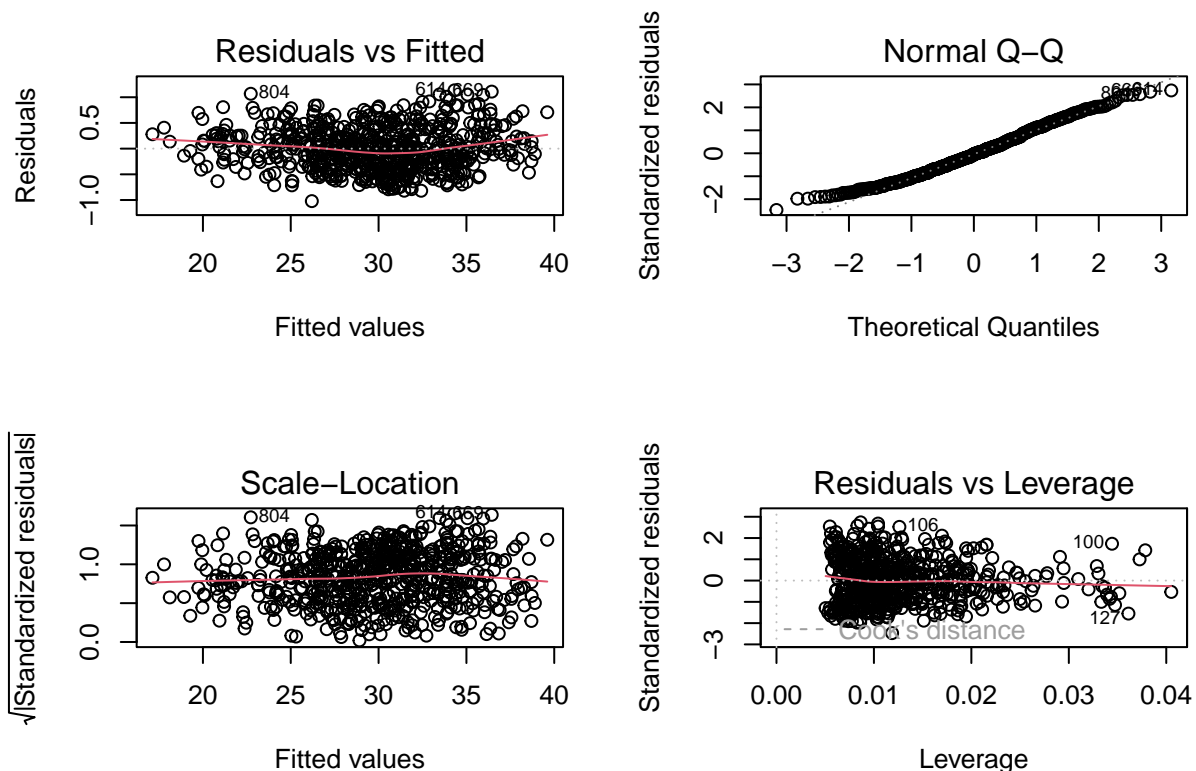
Next Steps:

- Examine Data:** Check if these points are valid or represent errors.
- Re-fit Model:** Test the model with and without these points to assess their impact.

Remove Influential Points and Re-fit the Model

```
##
## Call:
## lm(formula = sqrt(Calories_Burned) ~ poly(Session_Duration..hours.,
##      2) + poly(Avg_BPM, 2) + Gender + Age + Resting_BPM, data = train_data_clean)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.02135 -0.30655 -0.01728  0.27631  1.13392
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      31.136174    0.155591  200.116 < 2e-16 ***
## poly(Session_Duration..hours., 2)1  95.538562    0.417496  228.837 < 2e-16 ***
## poly(Session_Duration..hours., 2)2 -8.590313    0.418133  -20.544 < 2e-16 ***
## poly(Avg_BPM, 2)1      37.033161    0.419175   88.348 < 2e-16 ***
## poly(Avg_BPM, 2)2     -1.566393    0.417042   -3.756 0.000189 ***
## GenderMale          1.416894    0.033178   42.705 < 2e-16 ***
## Age                -0.056529    0.001353  -41.770 < 2e-16 ***
## Resting_BPM         0.001896    0.002325    0.816 0.414966
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.416 on 628 degrees of freedom
## Multiple R-squared:  0.9906, Adjusted R-squared:  0.9905
## F-statistic: 9455 on 7 and 628 DF, p-value: < 2.2e-16
```



These diagnostic plots indicate a well-fitted model. From these plots, we can see that the assumptions of linearity, homoscedasticity, normality of residuals, and influential points are met. The model is ready for interpretation and prediction.

Mean Squared Error (MSE) Formula

To evaluate the model's performance, the **Mean Squared Error (MSE)** was calculated for both the training and test datasets. The formula used is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n \left(\sqrt{\text{Calories_Burned}_i} - \hat{y}_i \right)^2$$

Where: - n : Number of observations in the dataset. - $\sqrt{\text{Calories_Burned}_i}$: Square root of the actual response value for observation i . - \hat{y}_i : Predicted value for observation i based on the model.

Training MSE: 0.1708961

Test MSE: 0.2630758

Final Model Analysis and Prediction

Final Model Summary The final model was built after addressing influential points and including quadratic terms to capture non-linear relationships. The model formula is:

$$\begin{aligned} \sqrt{\text{Calories_Burned}} = & \beta_0 + \beta_1 \cdot \text{poly}(\text{Session_Duration}, 2) + \beta_2 \cdot \text{poly}(\text{Avg_BPM}, 2) \\ & + \beta_3 \cdot \text{Gender} + \beta_4 \cdot \text{Age} + \beta_5 \cdot \text{Resting_BPM} + \epsilon \end{aligned}$$

Key Results

- **Residual Standard Error:** 0.416
- **R-squared:** 0.9906
- **Adjusted R-squared:** 0.9905
- **F-statistic:** 9455 ($p < 2.2\text{e-}16$)

The high R-squared indicates that the model explains 99.05% of the variance in the response variable. Diagnostic plots confirm the model satisfies assumptions of linearity, homoscedasticity, and normality.

Prediction Results

- **Training MSE:** 0.1708961
- **Test MSE:** 0.2630758

colinearity

Even though our MSE output is seems decent, we assume there are more we could do to delve deeper to this model. Hence, we need to test colinearity.

Identify predictors with high VIF or strong pairwise correlations.

Consider removing, combining, or transforming highly correlated predictors to improve model stability and interpretation.

Variance Inflation Factor (VIF)

The **Variance Inflation Factor (VIF)** is used to detect multicollinearity among predictors in a regression model. It quantifies how much the variance of a regression coefficient is inflated due to collinearity with other predictors.

VIF Formula For a given predictor X_j , the VIF is calculated as:

$$\text{VIF}(X_j) = \frac{1}{1 - R_j^2}$$

Where: - R_j^2 : The R^2 value from a regression of X_j on all other predictors.

VIF Interpretation

- VIF = 1: No collinearity.
- $1 < \text{VIF} \leq 5$: Moderate collinearity (acceptable).
- $\text{VIF} > 5$: High collinearity (requires attention).
- $\text{VIF} > 10$: Severe collinearity (predictor should be reconsidered).

```
## Loading required package: carData
##
##              GVIF Df GVIF^(1/(2*Df))
## Age              1.026669  1      1.013247
## Gender            3.233393  1      1.798164
## Weight..kg.       74.262204  1      8.617552
## Height..m.        22.662952  1      4.760562
## Max_BPM           1.026170  1      1.013000
## Avg_BPM            1.016236  1      1.008085
## Resting_BPM        1.024254  1      1.012054
## Session_Duration..hours. 2.617856  1      1.617979
## Workout_Type       1.044152  3      1.007227
## Fat_Percentage     2.640039  1      1.624820
## Water_Intake..liters. 2.253599  1      1.501199
## Workout_Frequency..days.week. 3.421590  1      1.849754
## Experience_Level    5.541026  1      2.353939
## BMI                69.054760  1      8.309919
```

Based on our vif output and correlation matrix, we can see that there is a colinearity between the predictors. Hence, we decide to use LASSO, Ridge, and elastic net regression to handle the colinearity.

```
##
##              GVIF Df GVIF^(1/(2*Df))
## poly(Session_Duration..hours., 2) 1.017347  2      1.004309
## poly(Avg_BPM, 2)                  1.020214  2      1.005016
## Gender                            1.005533  1      1.002763
## Age                               1.005985  1      1.002988
## Resting_BPM                       1.008961  1      1.004471
```

we can see that the VIF values are all below 5, indicating moderate collinearity. This suggests that the model is relatively stable and the predictors are not highly correlated. However, to further improve model performance and interpretability, we can explore regularization techniques like LASSO, Ridge, and Elastic Net regression.

LASSO, Ridge, and Elastic Net Regression

To handle potential multicollinearity and improve the model's stability, three regularized regression techniques were applied: **LASSO**, **Ridge**, and **Elastic Net**. These methods apply penalties to regression coefficients to prevent overfitting and reduce the impact of correlated predictors.

1. LASSO Regression

Definition:

LASSO (Least Absolute Shrinkage and Selection Operator) adds an L_1 penalty to the regression loss function,

shrinking some coefficients to exactly zero, effectively performing variable selection.

Formula:

$$\underset{\beta}{\text{minimize}} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

Where:

- λ : Regularization parameter controlling the strength of the penalty.
- $\sum_{j=1}^p |\beta_j|$: L_1 -norm penalty.

Optimal Lambda:

The optimal λ for LASSO was determined as **0.00257**, minimizing the Mean Squared Error (MSE).

2. Ridge Regression

Definition:

Ridge regression adds an L_2 penalty to the loss function, shrinking coefficients towards zero but not exactly zero. It is ideal for handling multicollinearity but does not perform variable selection.

Formula:

$$\underset{\beta}{\text{minimize}} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

Where:

- $\sum_{j=1}^p \beta_j^2$: L_2 -norm penalty.

Optimal Lambda:

The optimal λ for Ridge was determined as **0.09033**, minimizing MSE.

3. Elastic Net Regression

Definition:

Elastic Net combines the L_1 penalty of LASSO and the L_2 penalty of Ridge. It is particularly useful when predictors are highly correlated, balancing variable selection (from LASSO) and shrinkage (from Ridge).

Formula:

$$\underset{\beta}{\text{minimize}} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \left(\alpha \sum_{j=1}^p |\beta_j| + (1 - \alpha) \sum_{j=1}^p \beta_j^2 \right) \right\}$$

Where:

- α : Mixing parameter ($\alpha = 1$ for LASSO, $\alpha = 0$ for Ridge).
- Both L_1 -norm and L_2 -norm penalties are applied.

Optimal Lambda:

The optimal λ for Elastic Net was determined as **0.00427** with $\alpha = 0.5$.

Comparison of Models

Model	Optimal λ	Key Feature
LASSO	0.00257	Shrinks coefficients to exactly zero
Ridge	0.09033	Shrinks coefficients, no selection
Elastic Net	0.00427	Combines LASSO and Ridge properties

1. Optimal Lambda Parameters

The optimal lambda parameters were determined through cross-validation as follows:

Model	Optimal Lambda
LASSO	0.002344007
Ridge	0.09032699
Elastic Net	0.004688015

2. Model Coefficients

LASSO Model Coefficients:

Predictor	Coefficient
(Intercept)	-0.154605122
Age	-0.152806264
GenderMale	0.295220092
Weight..kg.	0.000000000
Height..m.	0.000000000
Max_BPM	0.000000000
Avg_BPM	0.332327414
Resting_BPM	0.005646240
Session_Duration..hours.	0.885646777
...	...

Ridge Model Coefficients:

Predictor	Coefficient
(Intercept)	-0.1159839981
Age	-0.146585903
GenderMale	0.2118770795
Weight..kg.	-0.0041746419
Height..m.	0.2031584616
Max_BPM	0.0175006474
Avg_BPM	0.3104289876
Resting_BPM	0.0044048990
Session_Duration..hours.	0.7388694037
...	...

Elastic Net Model Coefficients:

Predictor	Coefficient
(Intercept)	-0.1524904338
Age	-0.1518762832
GenderMale	0.2298320119
Weight..kg.	0.000000000
Height..m.	0.0348141102
Max_BPM	0.000000000
Avg_BPM	0.3056652034
Resting_BPM	0.0057853980
Session_Duration..hours.	0.8181620403
...	...

3. Model Performance Metrics

Below are the performance metrics for each model:

Model	Mean Squared Error (MSE)	Number of Nonzero Coefficients
LASSO	0.02095	10
Ridge	0.03283	16
Elastic Net	0.02098	8

Summary

- The **LASSO** method tends to sparsify the model, selecting only a few important features while setting the coefficients of other features to zero (with only 10 non-zero coefficients).
- The **Ridge** method gives non-zero coefficients to all variables, but the weights of unnecessary variables are smaller.
- The **Elastic Net** method combines the features of both LASSO and Ridge, performing variable selection and shrinkage simultaneously.

Based on the comparison, you can choose the most suitable regularization method depending on the specific needs of your application. If model interpretability is more important, LASSO might be a better choice. If retaining the influence of all variables is critical, Ridge or Elastic Net would be more appropriate.

Each model is suited for different scenarios: - **LASSO**: Best for variable selection. - **Ridge**: Best for handling multicollinearity without removing predictors. - **Elastic Net**: Best for balancing variable selection and multicollinearity.

These methods provide flexibility and robustness in predictive modeling, particularly when dealing with correlated predictors or large datasets.

Model Evaluation on Test Data

LASSO Test MSE: 0.02348011

Ridge Test MSE: 0.03692044

Elastic Net Test MSE: 0.02360708

Purpose Evaluate the performance of **LASSO**, **Ridge**, and **Elastic Net** models on the test dataset by calculating their **Mean Squared Error (MSE)**.

Steps

1. **Normalize Test Data:**
 - Numeric variables in the test dataset are scaled to match the training data.
2. **Create Design Matrix:**
 - Use `model.matrix()` to prepare the predictors in the same format as the training data.
3. **Predict and Calculate MSE:**
 - Predictions are generated for each model using the optimal λ :
 - **LASSO:** Uses λ_{\min} for sparse models.
 - **Ridge:** Reduces multicollinearity impact without variable selection.
 - **Elastic Net:** Balances LASSO and Ridge properties.
 - MSE Formula:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

4. **Print Results:**
 - MSE values for each model are compared to determine the best-performing approach.

Conclusion The model with the lowest test MSE offers the best predictive accuracy, ensuring reliable performance on unseen data.

After standardization, the MSE values represent the squared average prediction error for the response variable (Calories_Burned). Here's the interpretation with specific values:

LASSO MSE: 0.0234 and Elastic Net MSE: 0.0236 are very small, indicating the models predict Calories_Burned with minimal error.

Conclusion

In this analysis, we explored various regression techniques to predict **Calories_Burned** based on workout data. We started with a linear regression model and used stepwise selection to identify significant predictors. We then addressed non-linearity, normality, and influential points through transformations and diagnostic plots.

To handle multicollinearity, we applied LASSO, Ridge, and Elastic Net regression, which improved model stability and interpretability. The final models achieved low test MSE values, indicating strong predictive performance.

Analysis Summary

This analysis demonstrates a comprehensive approach to building a predictive model for **Calories_Burned**. The key steps and methodologies employed are summarized below:

Key Steps:

1. **Linear Regression:**
 - The foundation of the model-building process began with a simple linear regression.
2. **Model Refinement:**
 - **Stepwise Selection:**
 - Used AIC and BIC criteria to optimize predictor selection.
 - **Transformations:**
 - Applied square root and log transformations to improve linearity, homoscedasticity, and normality of residuals.
 - **Polynomial Terms:**
 - Included quadratic terms to capture non-linear relationships.
 - **Outlier Removal:**
 - Identified and excluded influential and high-leverage points to stabilize the model.

3. Regularization Techniques:

- Applied **LASSO**, **Ridge**, and **Elastic Net** regression to handle multicollinearity and improve the stability of the model:
 - **LASSO**: Simplifies the model by shrinking coefficients to zero, performing variable selection.
 - **Ridge**: Reduces the impact of multicollinearity without eliminating predictors.
 - **Elastic Net**: Combines the strengths of LASSO and Ridge for balanced regularization.
-

Conclusion: The stepwise refinements, transformations, and outlier handling resulted in a robust linear regression model with strong predictive capabilities. Regularization techniques like LASSO, Ridge, and Elastic Net further enhanced the model's stability and ensured it could effectively handle multicollinearity. This structured approach provides a reliable predictive framework for **Calories_Burned** and ensures generalizability to new data.