Analyzing Musical Feature Trends in Top Spotify Songs (2010–2019): A Statistical Overview

Jonah Sitorus

Dataset provided on Canvas (ProjectPhase2_ExcelDataFile_Spotify)
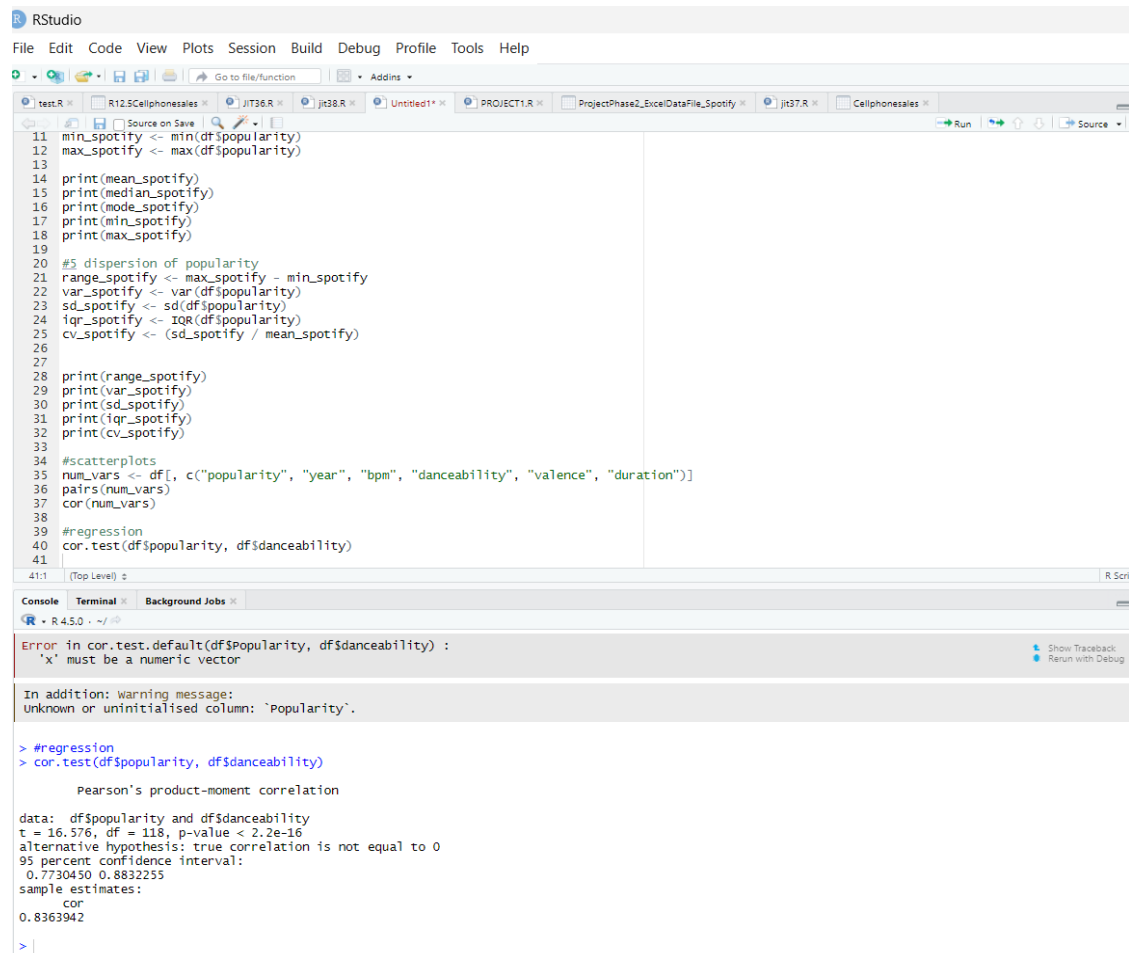
**Executive Summary**

In this project, I analyzed data on top Spotify songs from 2010-2019 to understand what makes a song popular. The dataset includes variables on musical features such as tempo (bpm), danceability, valence (mood), duration, year, and genre, along with a popularity score from 0 to 100. Understanding these relationships helps with identifying trends in listener preferences and better support playlist and marketing decisions.

Using R, I summarize the distribution of song popularity with descriptive statistics. I explore relationships between popularity and numerical variables using scatterplots and correlations. I built a multiple regression model with numerical predictors to explain variation in popularity and test model assumptions through residual analysis. I then extended the model by adding a categorical variable for the dance pop genre to see if this genre is associated with higher and lower popularity after adjusting for other musical attributes.

Overall, the results indicate that danceability, valence, and duration are statistically significant predictors of popularity, while year and tempo are not significant in this dataset. Residual models suggested that the assumptions were reasonably met and support the validity of this model. These findings provide an insight into the attributes that are in a top-charting Spotify song and highlights which feature is most popular and what is most geared towards listener preference.

**Descriptive Statistics for Popularity**

In this section, the descriptive statistics for the dependent variable, Popularity. The mean popularity was approximately 69.19222 with a median of 69.41. The minimum and maximum values were 59.90868 and 76.28763. The standard deviation was 3.216985 which indicates that the popularity scores were not spread out that much. Because the distribution is skewed to the right this means that the mean is greater than the median and the best measure of the center is the median because it's not heavily influenced by the high-value outliers in the right tail.

## Scatterplots and Correlation



h

The scatterplot of Popularity versus Danceability appears to show a roughly linear pattern and the direction is positive, meaning that as danceability increases, popularity also increases. The strength of the relationships appears strong based on how the points cluster around a trend and that the correlation level is high at .836 and there's no obvious outliers.



```
R RStudio

File   Edit   Code   View   Plots   Session   Build   Debug   Profile   Tools   Help

test.R ×   R12.5Cellphonesales ×   JIT36.R ×   jit38.R ×   Untitled1* ×   PROJECT1.R ×   ProjectPhase2_ExcelDataFile_Spotify ×   jit37.R ×   Cellphonesales ×

Source on Save                                                                                    → Run    ⇧ ⇩    Source ▾
 11  min_spotify <- min(df$popularity)
 12  max_spotify <- max(df$popularity)
 13
 14  print(mean_spotify)
 15  print(median_spotify)
 16  print(mode_spotify)
 17  print(min_spotify)
 18  print(max_spotify)
 19
 20  #5 dispersion of popularity
 21  range_spotify <- max_spotify - min_spotify
 22  var_spotify <- var(df$popularity)
 23  sd_spotify <- sd(df$popularity)
 24  iqr_spotify <- IQR(df$popularity)
 25  cv_spotify <- (sd_spotify / mean_spotify)
 26
 27
 28  print(range_spotify)
 29  print(var_spotify)
 30  print(sd_spotify)
 31  print(iqr_spotify)
 32  print(cv_spotify)
 33
 34  #scatterplots
 35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
 36  pairs(num_vars)
 37  cor(num_vars)
 38
 39  #regression
 40  cor.test(df$popularity, df$danceability)
 41
 41:1   (Top Level) ⇩                                                                              R Scri

Console   Terminal ×   Background Jobs ×

R ▾ R 4.5.0 · ~/

Error in cor.test.default(df$Popularity, df$danceability) :       Show Traceback
  'x' must be a numeric vector                                     Rerun with Debug

In addition: Warning message:
Unknown or uninitialised column: `Popularity`.

> #regression
> cor.test(df$popularity, df$danceability)

        Pearson's product-moment correlation

data:  df$popularity and df$danceability
t = 16.576, df = 118, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.7730450 0.8832255
sample estimates:
      cor
0.8363942

>
```

**Test significance of correlation (Popularity & Danceability)**
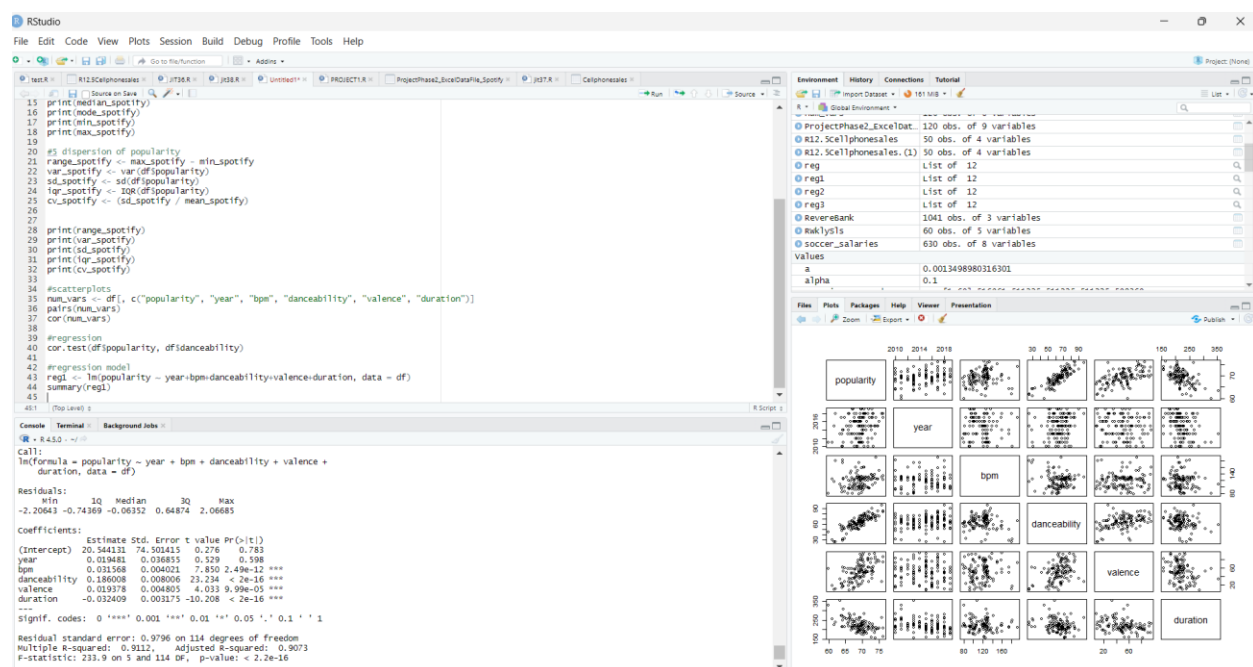
⬜ **H₀:** $\rho = 0$ (no linear correlation)

⬜ **Hₐ:** $\rho \neq 0$ (non-zero linear correlation)

The correlation between popularity and danceability is .836. Using cor.test(), we obtain a p-value of 2.2e^-16 which at the .05 significance level, we reject the null hypothesis of zero correlation. This indicates that there is a statistically significant linear association between popularity and danceability.

## Multicollinearity

For multicollinearity among the independent variables, I used the correlation matrix and pairs with correlations below 0.2 were classified as low, 0.2-0.8 as moderate and above 0.8 as extreme. In this dataset, popularity and year, popularity and bpm, popularity and duration had low multicollinearity, and popularity and valence had moderate multicollinearity and popularity and danceability and extreme multicollinearity.

## Multiple Regression Model



The r-squared for this model is 0.9112 meaning that ~91% of the variation in popularity is explained by the independent variables. This suggests the model has a strong fit to the data.

## Is R-Squared Statistically Significant?

☐ $H_0$: All slope coefficients = 0 (no linear relationship between Popularity and any predictor)

☐ $H_a$: At least one slope coefficient ≠ 0

The F-statistic is 233.9 with a p-value of 2.2e-16 and because this p-value is less than .05 we reject the null hypothesis. This means that the regression model is statistically significant in explaining the variation in popularity.

**Conclusion and interpretation**

$\beta 0$: if all of year, bpm, danceability, valence, and duration = 0, then estimated popularity = 20.544, not significant at the .05 level.
$\beta 1$: year does not impact popularity with statistical significance at the .05 level.
$\beta 2$: bpm does not impact popularity with statistical significance at the .05 level.

$\beta 3$: for a one unit increase in danceability, pop increases by .186, holding

all other variables fixed, significant at the .05 level.
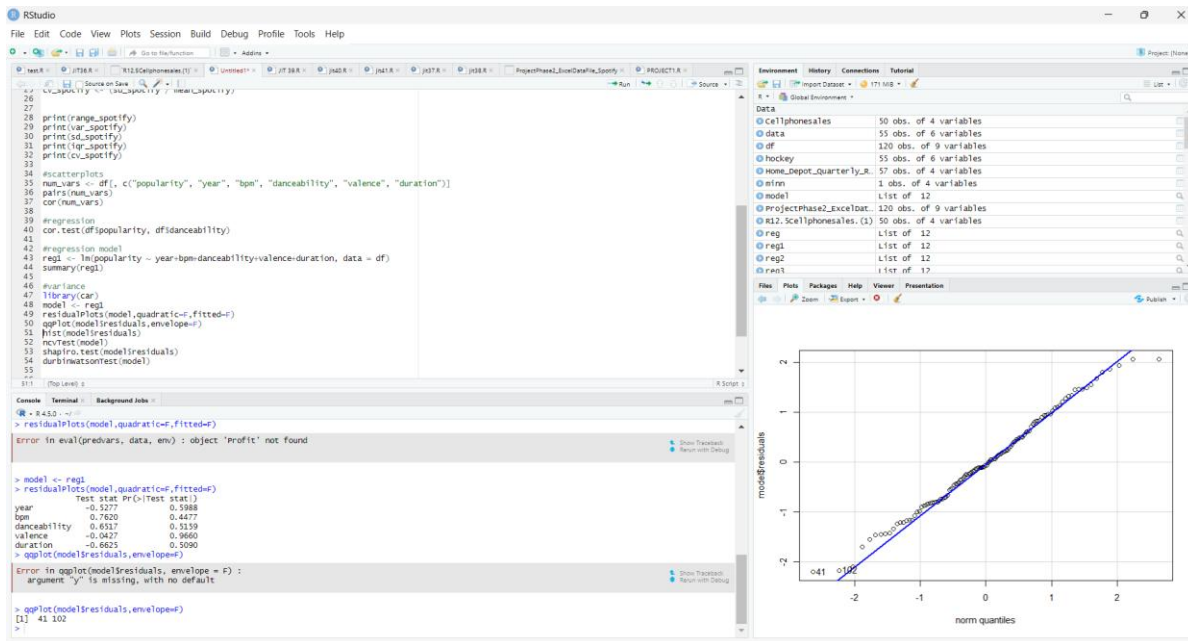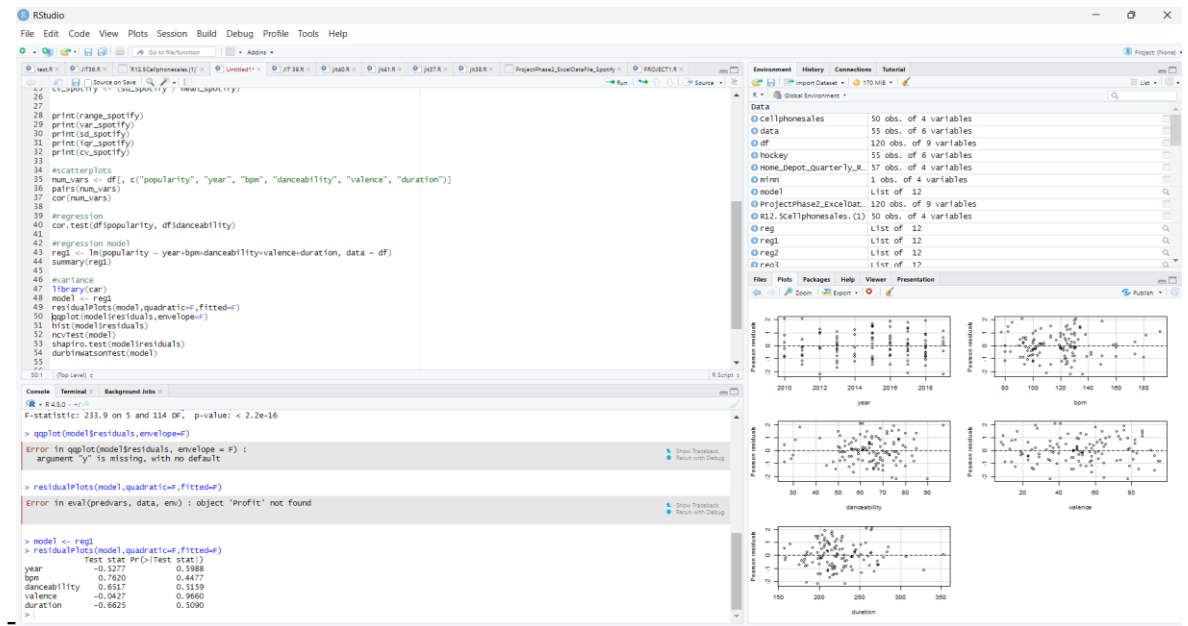
$\beta 4$: for a one unit increase in valence, pop increases by .193, holding

all other variables fixed, significant at the .05 level.

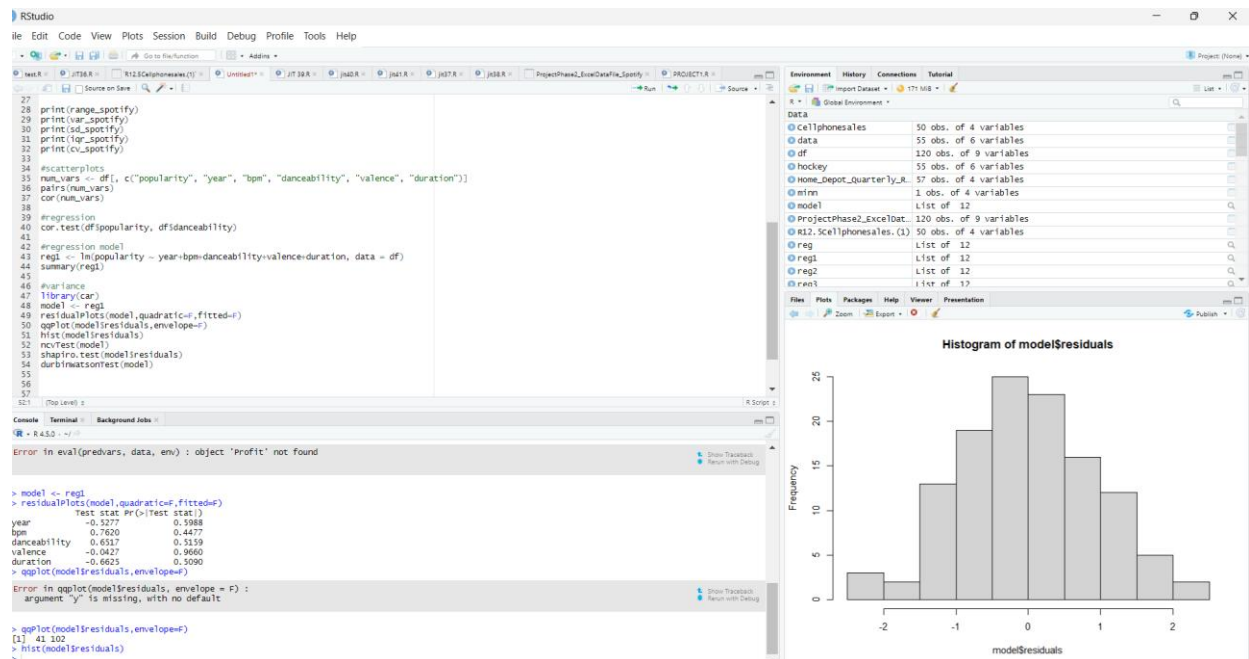$\beta 5$: for a one unit increase in duration, pop decreases by .324, holding

all other variables fixed, significant at the .05 level.

**Residual Analysis**

To assess whether the assumptions of simple linear regression were met, several statistical tests were used. These include graphical plots as well as formal tests for normality, homoscedasticity, and independence of residuals.

**Graphs**

**Screenshot 1 — Source editor:**

```
26
28  print(range_spotify)
29  print(var_spotify)
30  print(sd_spotify)
31  print(iqr_spotify)
32  print(cv_spotify)
33
34  #scatterplots
35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
36  pairs(num_vars)
37  cor(num_vars)
38
39  #regression
40  cor.test(df$popularity, df$danceability)
41
42  #regression model
43  reg1 <- lm(popularity ~ year+bpm+danceability+valence+duration, data = df)
44  summary(reg1)
45
46  #variance
47  library(car)
48  model <- reg1
49  residualPlots(model,quadratic=F,fitted=F)
50  qqplot(model$residuals,envelope=F)
51  hist(model$residuals)
52  ncvTest(model)
53  shapiro.test(model$residuals)
54  durbinwatsonTest(model)
55
```

**Console:**

```
F-statistic: 233.9 on 5 and 114 DF,  p-value: < 2.2e-16

> qqplot(model$residuals,envelope=F)

Error in qqplot(model$residuals, envelope = F) :
  argument "y" is missing, with no default

> residualPlots(model,quadratic=F,fitted=F)

Error in eval(predvars, data, env) : object 'Profit' not found

> model <- reg1
> residualPlots(model,quadratic=F,fitted=F)
            Test stat Pr(>|Test stat|)
year         -0.5277           0.5988
bpm           0.7620           0.4477
danceability  0.6517           0.5159
valence      -0.0427           0.9660
duration     -0.6625           0.5090
>
```

**Screenshot 2 — Source editor:**

```
26
27
28  print(range_spotify)
29  print(var_spotify)
30  print(sd_spotify)
31  print(iqr_spotify)
32  print(cv_spotify)
33
34  #scatterplots
35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
36  pairs(num_vars)
37  cor(num_vars)
38
39  #regression
40  cor.test(df$popularity, df$danceability)
41
42  #regression model
43  reg1 <- lm(popularity ~ year+bpm+danceability+valence-duration, data = df)
44  summary(reg1)
45
46  #variance
47  library(car)
48  model <- reg1
49  residualPlots(model,quadratic=F,fitted=F)
50  qqplot(model$residuals,envelope=F)
51  hist(model$residuals)
52  ncvTest(model)
53  shapiro.test(model$residuals)
54  durbinwatsonTest(model)
55
```

**Console:**

```
> residualPlots(model,quadratic=F,fitted=F)

Error in eval(predvars, data, env) : object 'Profit' not found

> model <- reg1
> residualPlots(model,quadratic=F,fitted=F)
            Test stat Pr(>|Test stat|)
year         -0.5277           0.5988
bpm           0.7620           0.4477
danceability  0.6517           0.5159
valence      -0.0427           0.9660
duration     -0.6625           0.5090
> qqplot(model$residuals,envelope=F)

Error in qqplot(model$residuals, envelope = F) :
  argument "y" is missing, with no default

> qqplot(model$residuals,envelope=F)
[1] 41 102
>
```

The residual plot reveals a clear pattern in year and its residual value and it slowly gets random in the other independent variables. This could mean that the linear assumption is met and that there's evidence for homoscedasticity if the plot shows a clear pattern. For the Q-Q plot, there is a strong linear pattern in the graph which could indicate normal linearity. The histogram is unimodal and approximately symmetric, indicating a normal distribution of the data.

**Breusch-Pagan Test**

File   Edit   Code   View   Plots   Session   Build   Debug   Profile   Tools   Help

test.R   JIT36.R   `R12.5Cellphonesales.(1)`   Untitled1*   JIT 39.R   jit40.R   jit41.R   jit37.R   jit38.R   ProjectPhase2_ExcelDataFile_Spotify   PROJECT1.R

```
27
28  print(range_spotify)
29  print(var_spotify)
30  print(sd_spotify)
31  print(iqr_spotify)
32  print(cv_spotify)
33
34  #scatterplots
35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
36  pairs(num_vars)
37  cor(num_vars)
38
39  #regression
40  cor.test(df$popularity, df$danceability)
41
42  #regression model
43  reg1 <- lm(popularity ~ year+bpm+danceability+valence+duration, data = df)
44  summary(reg1)
45
46  #variance
47  library(car)
48  model <- reg1
49  residualPlots(model,quadratic=F,fitted=F)
50  qqPlot(model$residuals,envelope=F)
51  hist(model$residuals)
52  ncvTest(model)
53  shapiro.test(model$residuals)
54  durbinwatsonTest(model)
55
56
57
```

53:1   (Top Level)                                                                                                R Scrip

Console   Terminal   Background Jobs

R · R 4.5.0 · ~/
```
valence         -0.0427          0.9660
duration        -0.6625          0.5090
> qqplot(model$residuals,envelope=F)

Error in qqplot(model$residuals, envelope = F) :
  argument "y" is missing, with no default


> qqPlot(model$residuals,envelope=F)
[1]  41 102
> hist(model$residuals)
> residualPlots(model,quadratic=F,fitted=F)
            Test stat Pr(>|Test stat|)
year          -0.5277           0.5988
bpm            0.7620           0.4477
danceability   0.6517           0.5159
valence       -0.0427           0.9660
duration      -0.6625           0.5090
> ncvTest(model)
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 2.615788, Df = 1, p = 0.1058
>
```

Show Traceback
Rerun with Debug

Interpretation:

- If p > 0.05 → Homoscedasticity holds (constant variance).
- If p < 0.05 → Heteroscedasticity present.
- The Breusch–Pagan test resulted in p > 0.05, suggesting the variance of residuals is constant, and we fail to reject that homoscedasticity is present.

**Shapiro-Wilk Test**

```
32  print(cv_spotify)
33
34  #scatterplots
35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
36  pairs(num_vars)
37  cor(num_vars)
38
39  #regression
40  cor.test(df$popularity, df$danceability)
41
42  #regression model
43  reg1 <- lm(popularity ~ year+bpm+danceability+valence+duration, data = df)
44  summary(reg1)
45
46  #variance
47  library(car)
48  model <- reg1
49  residualPlots(model,quadratic=F,fitted=F)
50  qqPlot(model$residuals,envelope=F)
51  hist(model$residuals)
52  ncvTest(model)
53  shapiro.test(model$residuals)
54  durbinwatsonTest(model)
55
56
57
58  #new variable
59  df$dancepop <- ifelse(df$genre == "dance pop", 1, 0)
60  reg2 <- lm(popularity ~ year + bpm + danceability + valence + duration + dancepop, data = df)
61  summary(reg2)
62
```

```
> qqPlot(model$residuals,envelope=F)
[1]  41 102
> hist(model$residuals)
> residualPlots(model,quadratic=F,fitted=F)
             Test stat Pr(>|Test stat|)
year           -0.5277           0.5988
bpm             0.7620           0.4477
danceability    0.6517           0.5159
valence        -0.0427           0.9660
duration       -0.6625           0.5090
> ncvTest(model)
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 2.615788, Df = 1, p = 0.1058
> shapiro.test(model$residuals)

        Shapiro-wilk normality test

data:  model$residuals
W = 0.99059, p-value = 0.588

>
```

Interpretation:

- If p > 0.05 → Fail to reject normality

- If p < 0.05 → Evidence of non-normality.

- The Shapiro–Wilk test yielded p > 0.05, indicating that we fail to reject normality

- The result aligns with what is seen visually in the Q–Q plot and histogram.

**Durbin–Watson Test**

```
 32  print(cv_spotify)
 33
 34  #scatterplots
 35  num_vars <- df[, c("popularity", "year", "bpm", "danceability", "valence", "duration")]
 36  pairs(num_vars)
 37  cor(num_vars)
 38
 39  #regression
 40  cor.test(df$popularity, df$danceability)
 41
 42  #regression model
 43  reg1 <- lm(popularity ~ year+bpm+danceability+valence+duration, data = df)
 44  summary(reg1)
 45
 46  #variance
 47  library(car)
 48  model <- reg1
 49  residualPlots(model,quadratic=F,fitted=F)
 50  qqPlot(model$residuals,envelope=F)
 51  hist(model$residuals)
 52  ncvTest(model)
 53  shapiro.test(model$residuals)
 54  durbinwatsonTest(model)
 55
 56
 57
 58  #new variable
 59  df$dancepop <- ifelse(df$genre == "dance pop", 1, 0)
 60  reg2 <- lm(popularity ~ year + bpm + danceability + valence + duration + dancepop, data = df)
 61  summary(reg2)
 62
```

```
> qqPlot(model$residuals,envelope=F)
[1]  41 102
> hist(model$residuals)
> residualPlots(model,quadratic=F,fitted=F)
             Test stat Pr(>|Test stat|)
year           -0.5277           0.5988
bpm             0.7620           0.4477
danceability    0.6517           0.5159
valence        -0.0427           0.9660
duration       -0.6625           0.5090
> ncvTest(model)
Non-constant Variance Score Test
Variance formula: ~ fitted.values
Chisquare = 2.615788, Df = 1, p = 0.1058
> shapiro.test(model$residuals)

	Shapiro-Wilk normality test

data:  model$residuals
W = 0.99059, p-value = 0.588

>
```
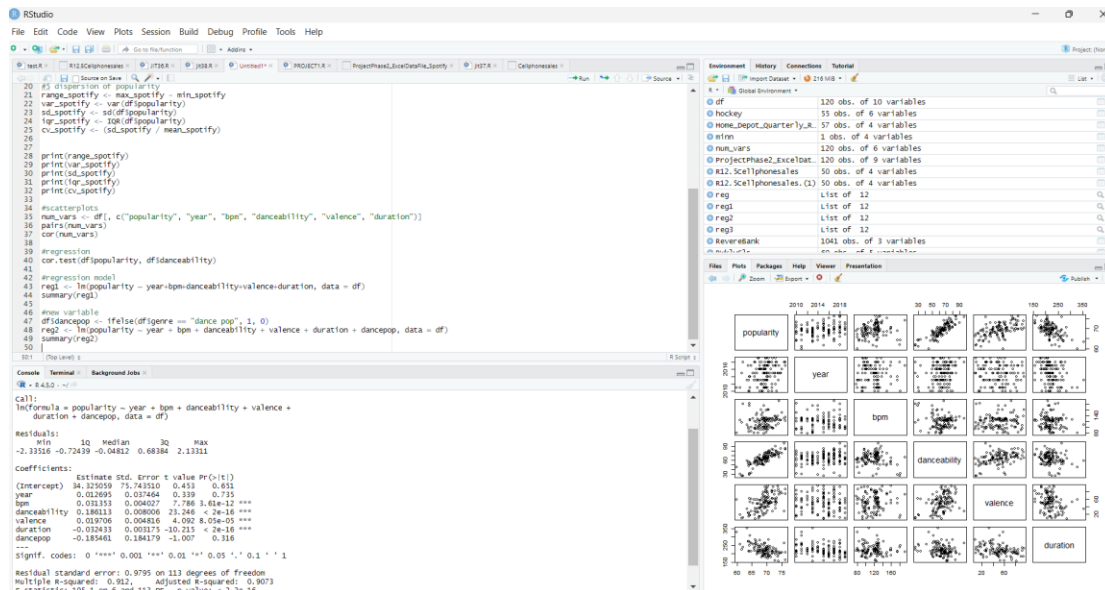
**4**

**Interpretation:**

- p > 0.05 → No significant autocorrelation (independent errors)

- p < 0.05 → Evidence of autocorrelation

- The Durbin–Watson test produced p > 0.05, indicating that there is no significant autocorrelation

- This suggests the independence assumption is met.

All diagnostic tools and statistical tests indicate that the regression assumptions are satisfied and so the regression model provides a reliable and unbiased fit for the data.

## Extended Model Including Genre

The adjusted r-squared for the original model was .9073 and when the dancepop variable was added the adjusted r-squared became .9073. Because the adjusted r-squared was the same it didn't improve the model's explanatory power.

$\beta6$: The coefficient on dance pop is -.1854 with a p-value of .316. This means that with all the variables held constant, songs classified as dance pop are predicted to have .1854 less popularity points than non-dance pop songs. The p-value is greater than 0.05, this is statistically significant at the .05 level.