# Objective

The objective of our analysis of the Spotify songs it to find a solid linear model that can be used to predict the danceability of a song based on certain variables on the dataset.

# Dataset

- 42305 songs on Spotify with **twenty-two** variables: Danceability, Energy, Key, Loudness, Mode, Speechiness, Acousticness, Intrumentalness, Liveness, Valence, Tempo, Type, ID, URI, Track_href, Analysis-url, Duration_ms, Time_signature, Genre, Song_name, Unnamed..O and Title.

- 42305 songs on Spotify with eleven variables: Danceability, Energy, key, Loudness, Speechiness, Acousticness, Intrumentalness, Liveness, Valence, Tempo, Duration_ms

# What is Danceability?

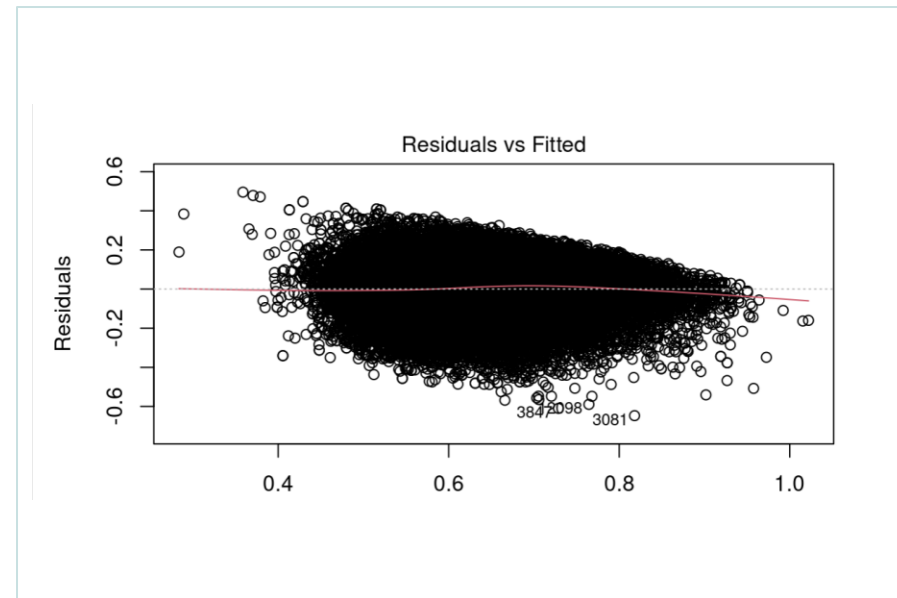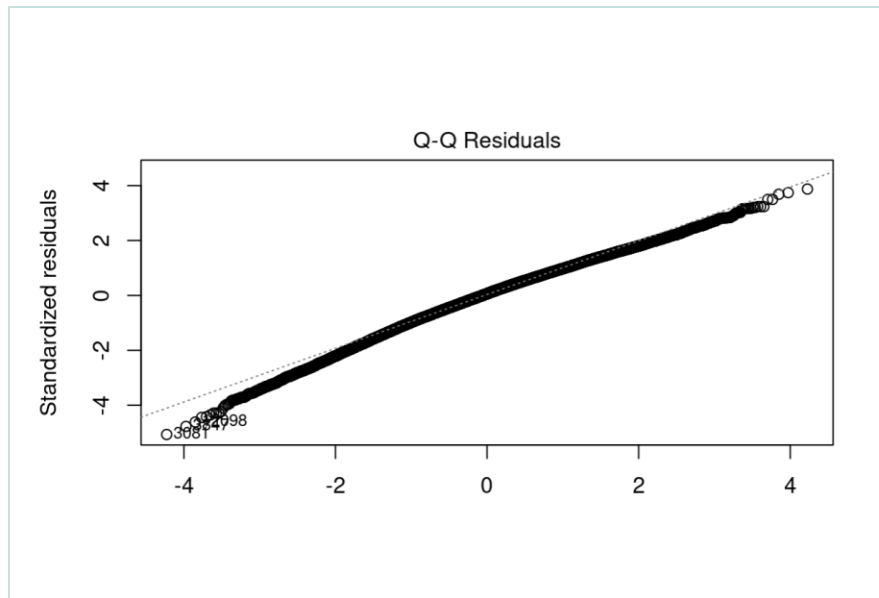Danceability describes how suitable a song is for dancing.

A numerical variable which is from 0-1 with 0 being least danceable and 1 being most danceable.
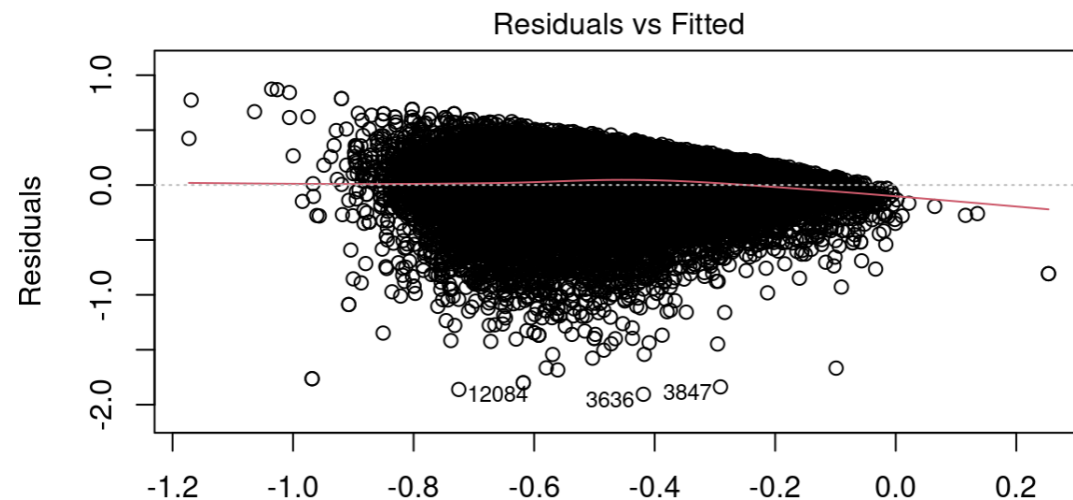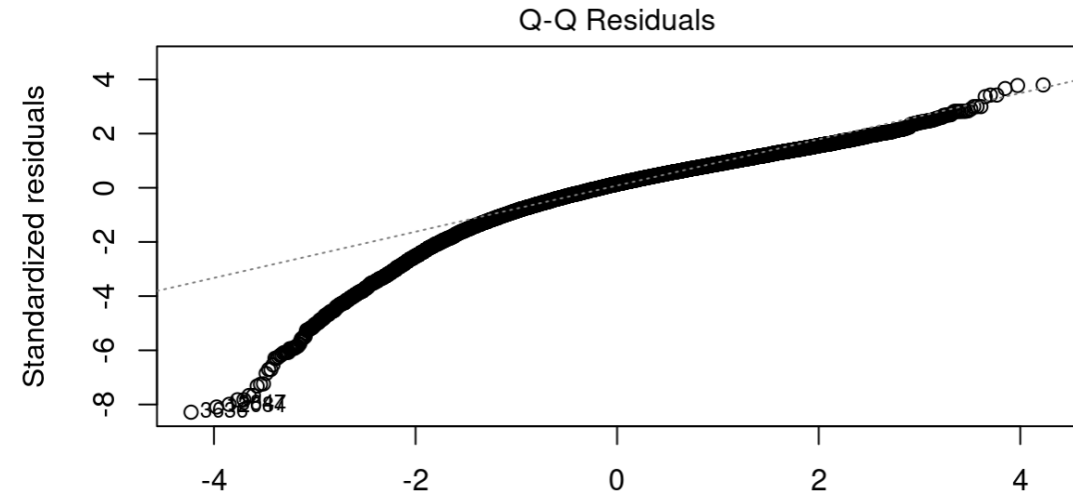
Mean = 0.639

# Finding The Best Model

- Ran stepwise regression using regsubsets to determine best model, maximizing for adjusted R-squared.

- Stepwise regression determined that the maximum model with all ten predictors in the data set had the highest adjusted R-squared.
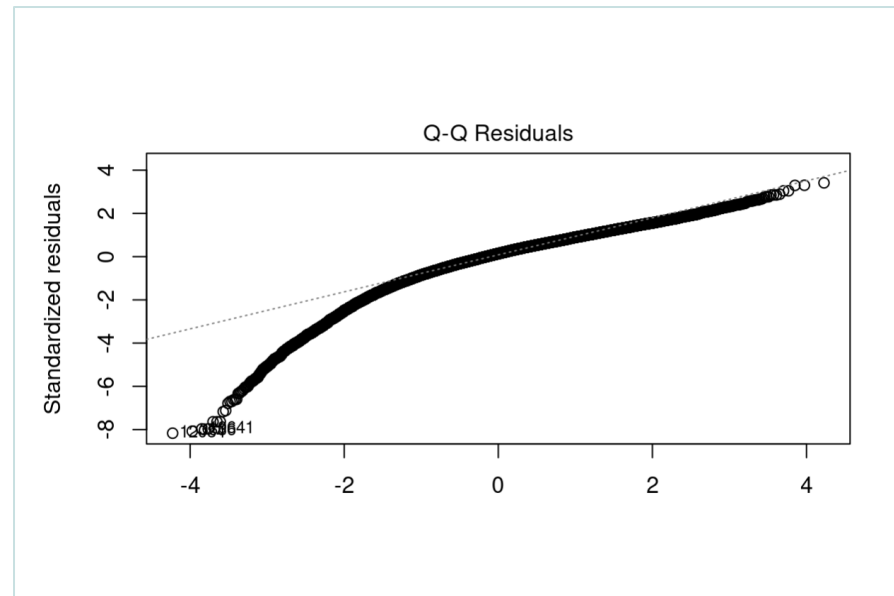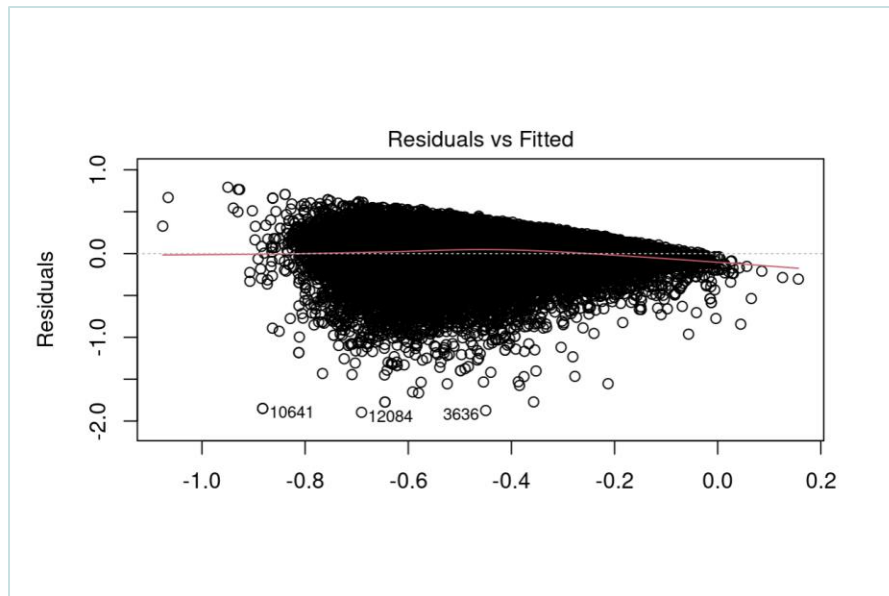
- Adjusted R-squared: .3367

# Assumptions for Maximized Model

# Adding Interaction Terms

- Added interaction terms between tempo and energy, as well as instrumentalness and loudness.

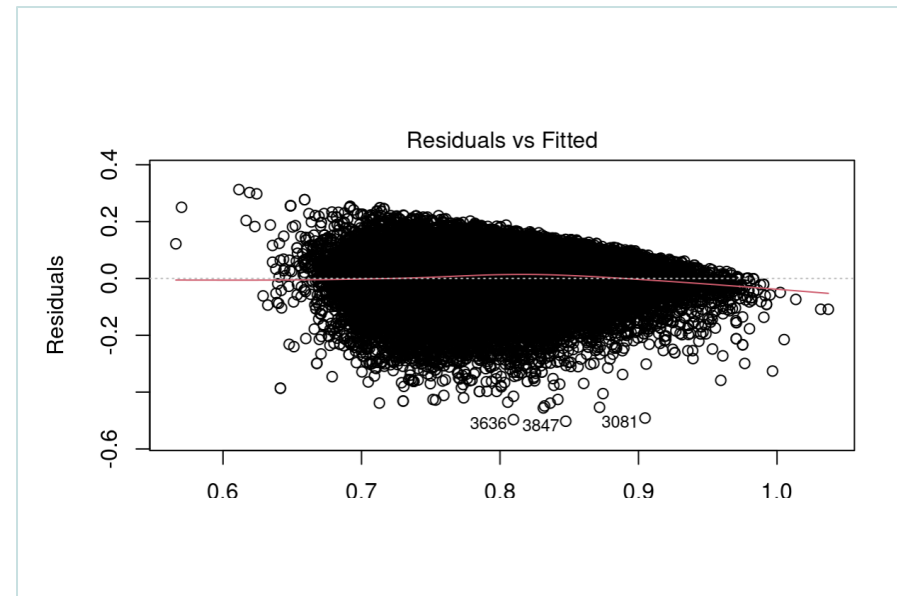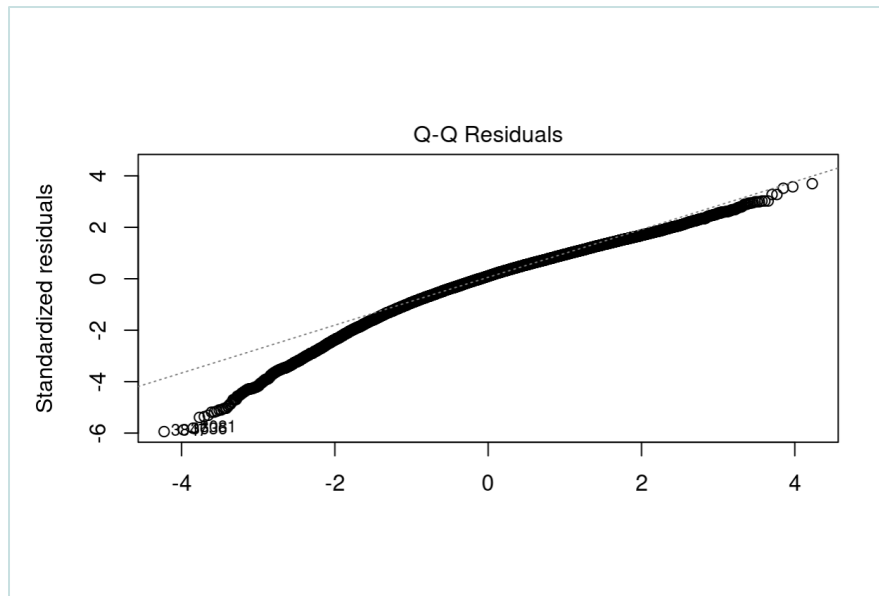- Adjusted R-squared: .3340, still statistically significant.

Adjusted R-squared: .2932

# Log Transformation

Adjusted R-Squared:.3173

# Square Root Transformation

# Key Takeaways

- The full model ended up being the most accurate model to predict the danceability of a song.

- Only ~33% of the dataset's variation in song danceability could be explained by the linear model

- Interaction terms were not significant to our model.

- Log transformation did not improve model.

- Dataset may not be suited to linear regression, further exploration through non-linear methods may be appropriate.

Questions?