

A Bayesian Analysis of Spotify Data

Nathaniel Maxwell, Jessie Bierschenk

01 May, 2021

Introduction

Music-making is often thought of as an artform—a subjective expression that falls into a specific “genre” according to its musical attributes. Beginning in the 1960s, pop music had been dominated by the verse-chorus form where “the verse sets the scene, the pre-chorus builds tension, and the chorus reaches a climax,” with the cycle predictably repeating itself [NY TIMES]. This musical formula dominated the industry; in fact, “music theorist Jay Summach has found that by the end of the 1960s, 42 percent of hit songs used verse-chorus form. By the end of the 1980s, that figure had doubled to 84 percent” [NY TIMES]. With the advent of the 21st Century, however, the digitization of music production paired with the introduction of streaming platforms has warped the fundamental structure of songs. On popular media platforms, only snapshots of songs reach the ears of the public: five-second memes, 15-second TikToks, or 30-second ads. The limitless access to songs on streaming platforms has changed the landscape of song-making—“the gist of it: songwriters now get to the good stuff sooner” [Economist]. This phenomenon exists as increasing accessibility of songs results in decreasing revenue for artists. “Artists are paid per play—provided the listener stays tuned for at least 30 seconds. Each stream earns a tiny fraction of a cent. And just 13% of that goes to the songwriter, says David Israelite of the National Music Publishers Association” [Economist]. In turn, for an artist to make a decent living, their songs need millions of plays. For many musicians, the art of composing/performing/marketing a new song is an arduous process. Even after all the work has been completed and a song is ready to be played to the public, the biggest uncertainty still awaits: How will the song be received? Will it become a hit? Will it be a song that everyone skips over, or never becomes popular? The purpose of this analysis is to investigate which characteristics of a song (such as tempo, duration, mode, acousticalness, etc.) would make it more “likeable,” less likely to be skipped, or more popular. Of course, music taste is a very subjective matter, and thus, there will be quite a bit of uncertainty around any variables that are deemed important/unimportant. What one person likes; another person may dislike. Therefore, looking at such musical characteristics through a Bayesian lens will help to quantify the uncertainty surrounding any of our findings. Through this analysis we hope to provide some conclusions that an aspiring musician (or even a well-established musician) can have at their disposal when creating new music. These findings beg the question: so what are the features that make a song popular or appealing to a listener? Answers that would be valuable for any musician seeking success in today’s music industry.

Pre-Analysis

Data

Two datasets were utilized during this analysis.

1. The first dataset consists of 83,939 observations on Spotify of whether or not a track was skipped by users. In total, 65,417 different tracks were included in the dataset. Each track has the following characteristics:

- (a) Release Year (Year the song was released)
- (b) Duration (length of song in seconds)
- (c) US Popularity Estimate (A popularity rating of song, on a scale 1-100)
- (d) Acousticness (A confidence measure from 0-1 on whether the track is acoustic, where values near 1 represent high confidence that the track is acoustic)
- (e) Beat Strength (The strength of the beat from 0-1, where 1 represents a very strong sense of beat)
- (f) Bounciness (A rating of the bounciness from 0-1, where 1 represents a strong sense of bounciness)
- (g) Danceability (A rating from 0-1 of how suitable the track is for dancing, where values near 1 represent high suitability)
- (h) Energy (A rating from 0-1 representing a perceptual measure of intensity and activity, where values near 1 represent high energy)
- (i) Instrumentalness (A rating from 0-1 that predicts whether a track has no vocals, where values close to 1 represent high confidence that there are no vocals)
- (j) Mode (Predicts whether or not a song is major or minor)
- (k) Speechiness (A rating from 0-1 that detects the presence of spoken words in a track, with values near 1 representing an exclusively speech-like track)
- (l) Tempo (The estimated tempo of the track in Beats Per Minute (BPM))
- (m) Valence (A rating from 0-1 that represents the positivity of the song, with 1 representing high positivity)
- (n) Skipped (Denotes whether or not that particular track was skipped or played the entire way through)

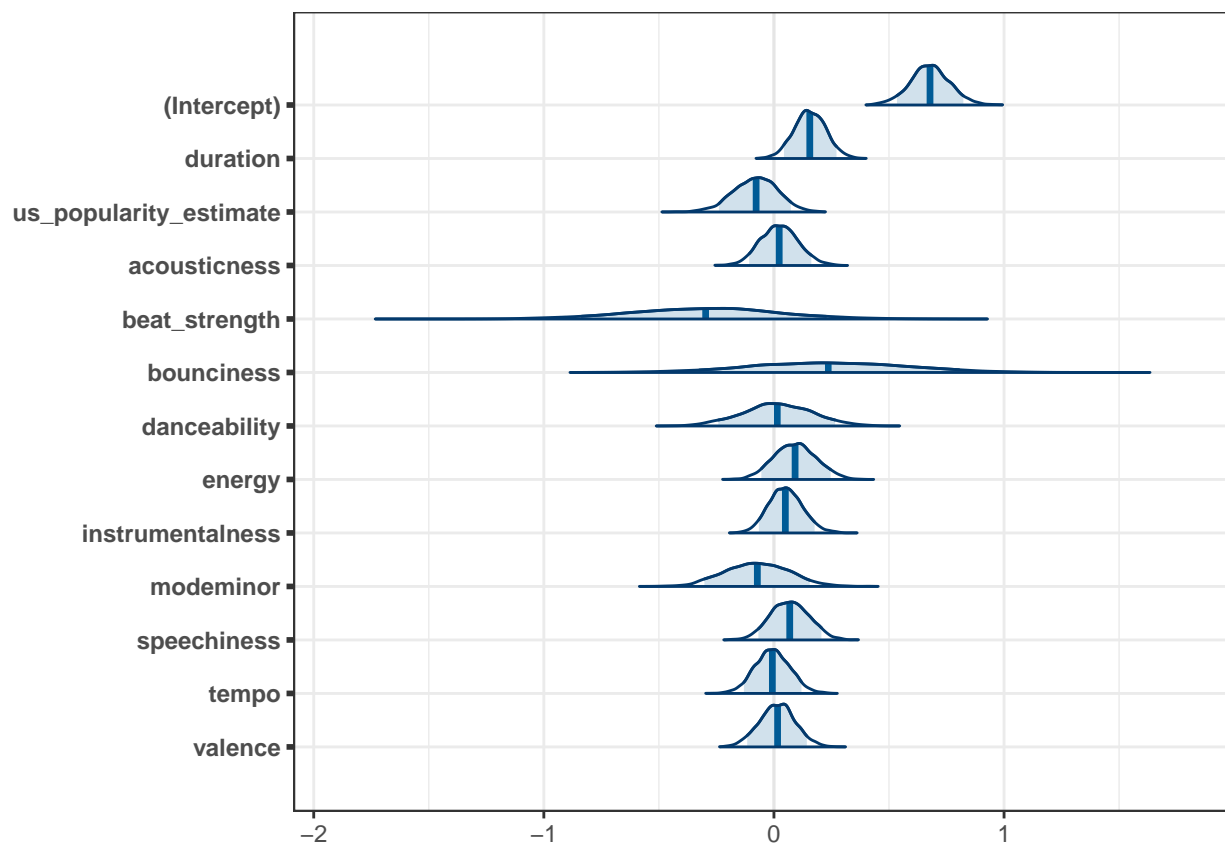
Note: in order to try to obtain tracks most representative of new music, only the following tracks were kept:

- (a) Tracks from 2010-present
 - (b) Tracks with a speechiness value ≤ 0.4 (filters out tracks that are mostly spoken, such as podcasts and ebooks)
 - (c) Tracks with an instrumentalness value ≤ 0.6 (filters out tracks that contain no vocals)
 - (d) Tracks with a duration ≤ 360 seconds (given that the average new song is 3-5 minutes, a cutoff of 6 minutes seemed appropriate)
2. The second dataset consisted of 2017 songs compiled by a single person, where a portion of the songs are songs that he likes, and the other portion are songs that he dislikes. This dataset includes similar variables as the first dataset, including:
- (a) Acousticness
 - (b) Danceability
 - (c) Duration
 - (d) Energy
 - (e) Instrumentalness
 - (f) Key (The particular grouping of chords and notes in a song)
 - (g) Liveness (rating from 0-1 of whether the track was performed live, with 1 representing high confidence the track was performed live)
 - (h) Loudness (Overall loudness of the track in decibels (dB))
 - (i) Mode
 - (j) Speechiness
 - (k) Tempo
 - (l) Time Signature (The way in which beats of the song are organized)
 - (m) Valence

Model Selection

For the first dataset, we wanted to estimate the values of the coefficients β for each of the variables to find out how they impact whether or not a track is skipped. We are assuming little knowledge about each variable's effect, so we propose a weakly informative prior for $[\beta]$: Using recommendations from Gelman, Jakulin, Pittau, and Su, we use a $\text{cauchy}(0, 2.5)$ prior for each scaled variable (we scaled the variables). Our response variable, \mathbf{y} , will follow a logistic regression model, where 1 means the track was skipped. This is equivalent to the Bernoulli distribution $\mathbf{y}|\theta \sim \text{Bern}(\theta)$. We will use the logit link, where $\text{logit}(\theta) = \eta$, and $\eta = \mathbf{x}^T \beta$, where \mathbf{x} is the covariate space for \mathbf{Y} . Using the `rstanarm` package, Rstudio will compute the posterior and draw MCMC samples from the posterior distribution $[\beta|\mathbf{Y}, \mathbf{X}]$.

Posterior Estimates



##	5%	95%
## (Intercept)	0.535	0.823
## duration	0.043	0.272
## us_popularity_estimate	-0.239	0.073
## acousticness	-0.108	0.162
## beat_strength	-0.836	0.227
## bounciness	-0.324	0.804
## danceability	-0.242	0.266
## energy	-0.055	0.246
## instrumentalness	-0.065	0.177
## modeminor	-0.302	0.156

```
## speechiness          -0.068 0.206
## tempo                -0.130 0.120
## valence              -0.116 0.143

##
## Computed from 4000 by 1000 log-likelihood matrix
##
##           Estimate    SE
## elpd_loo   -654.9 10.2
## p_loo       13.2  0.8
## looic      1309.8 20.3
## -----
## Monte Carlo SE of elpd_loo is 0.1.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.

##
## Computed from 4000 by 1000 log-likelihood matrix
##
##           Estimate    SE
## elpd_loo   -648.5  9.3
## p_loo        1.0  0.0
## looic      1296.9 18.6
## -----
## Monte Carlo SE of elpd_loo is 0.0.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.

##           elpd_diff se_diff
## posterior0  0.0         0.0
## posterior1 -6.5         3.7
```

After running the `rstanarm` function and including all of the variables, we see that there is only variable whose 90% confidence interval does not include 0. That variable is duration, and furthermore, when calculating the 'leave-one-out' cross-validation information criterion (looic), we see that this model actually has a *higher* value than the looic of a baseline model with no predictors. In other words, our model is worse at predicting whether or not a song is skipped than if someone randomly guessed! Therefore, we will drop all variables that were not deemed significant at a 90% confidence interval (included 0 in their posterior interval), and rerun the model. In this case, 'duration' is the only variable remaining.

```
posterior2 <- stan_glm(skipped ~ duration, data = Track_features_a,
                      family = binomial(link = "logit"),
                      prior = cauchy(0,2.5), prior_intercept = cauchy(0,2.5),
                      seed = seed,
                      refresh = 0)
(loo2 <- loo(posterior2, save_psis = TRUE))
```

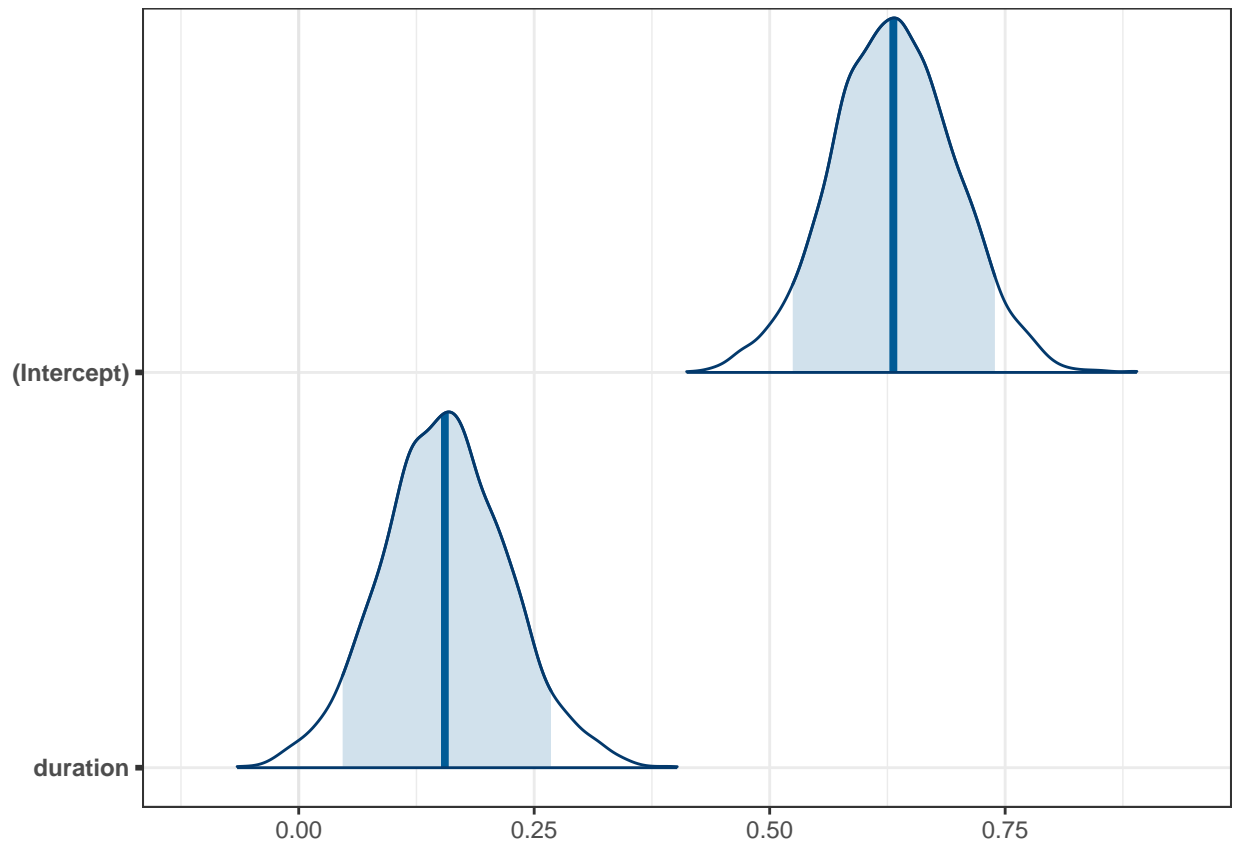
```
##
## Computed from 4000 by 1000 log-likelihood matrix
##
```

```
##           Estimate   SE
## elpd_loo   -646.7  9.6
## p_loo       2.0  0.1
## looic      1293.4 19.3
## -----
## Monte Carlo SE of elpd_loo is 0.0.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.
```

```
rstanarm::loo_compare(loo0, loo2)
```

```
##           elpd_diff se_diff
## posterior2  0.0      0.0
## posterior0 -1.8      2.4
```

```
mcmc_areas(as.matrix(posterior2), prob = 0.90, prob_outer = 1)
```



```
round(posterior_interval(posterior2, prob = 0.90), 3)
```

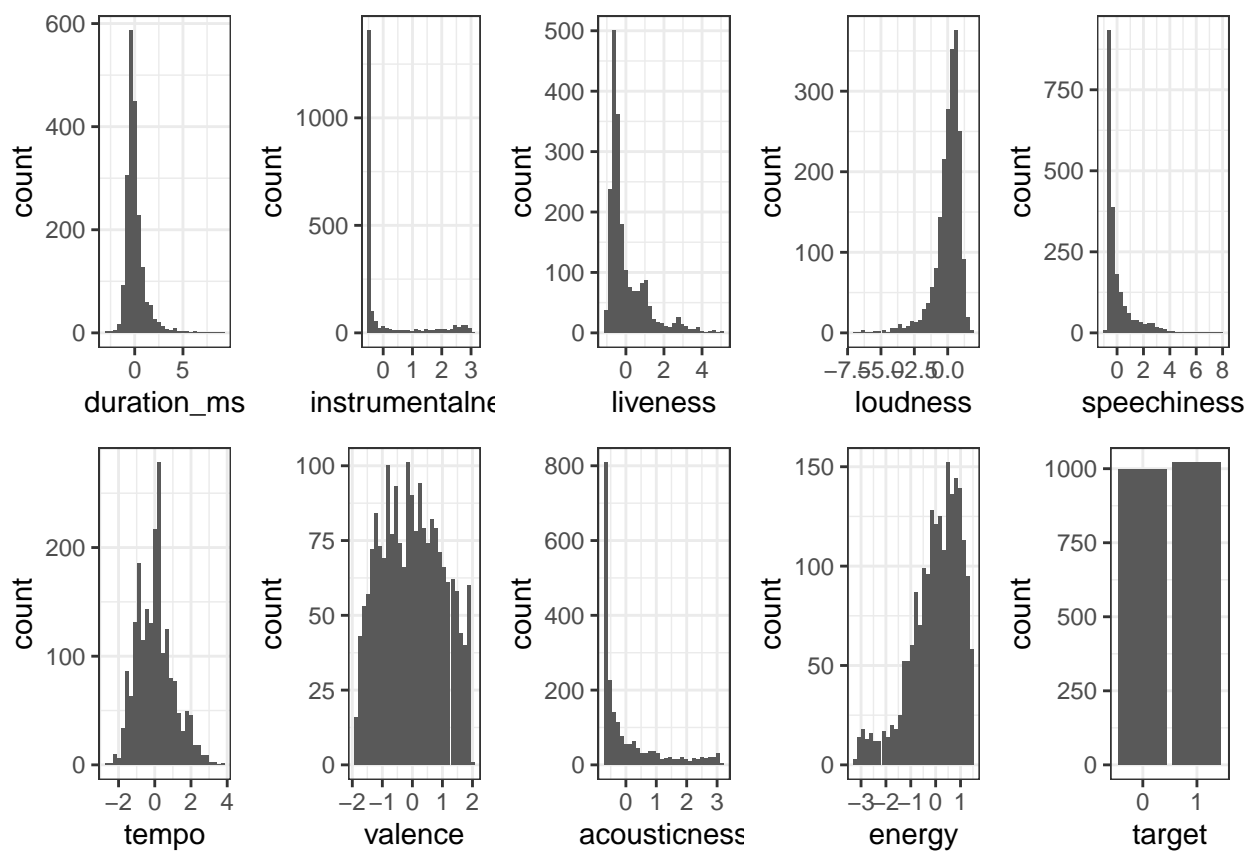
```
##           5%   95%
## (Intercept) 0.524 0.739
## duration    0.047 0.268
```

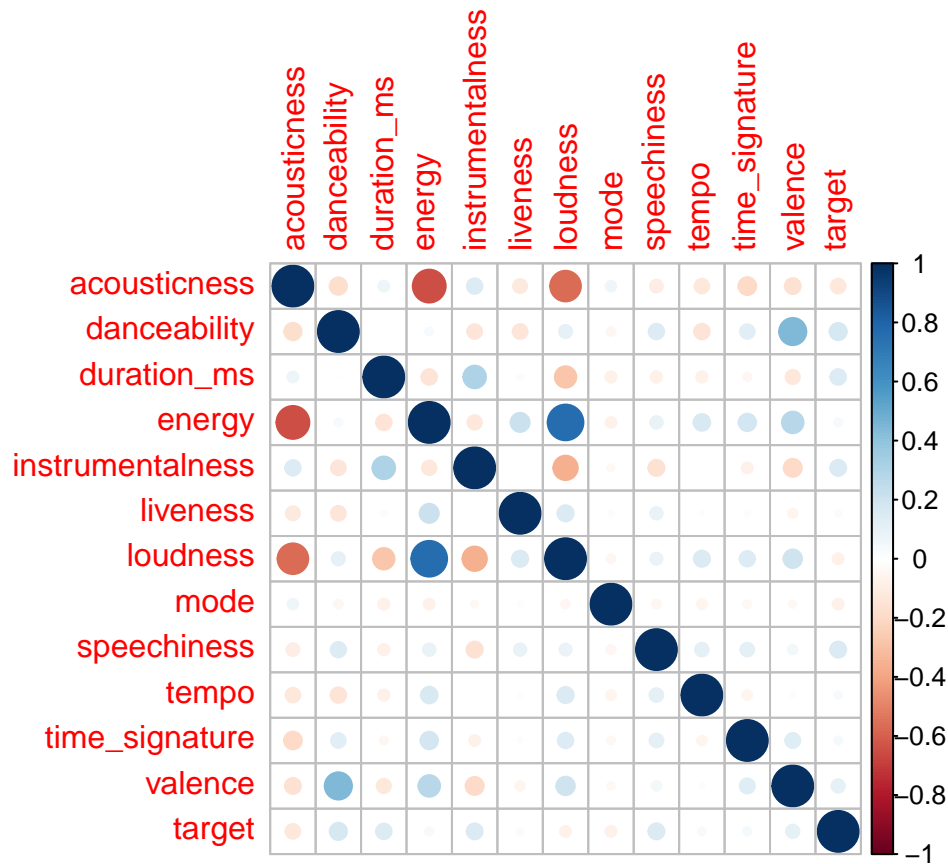
This model proved to be better, but not by much. Furthermore, the p

```
## [1] 0.65
```

```
## [1] 0.65
```

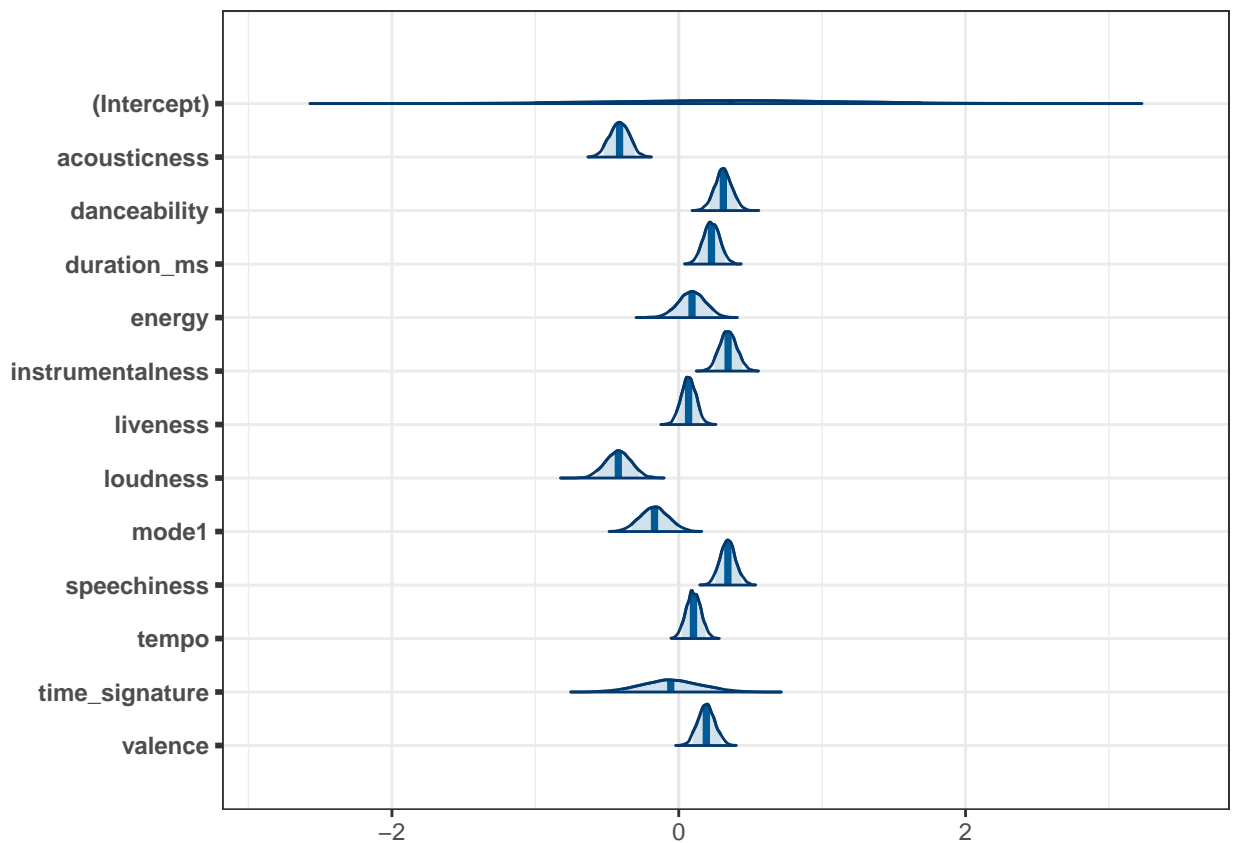
```
###New Data
```





```
##
## Model Info:
## function:      stan_glm
## family:        binomial [logit]
## formula:       target ~ .
## algorithm:     sampling
## sample:        4000 (posterior sample size)
## priors:        see help('prior_summary')
## observations:  2017
## predictors:    13
##
## Estimates:
##           mean    sd  10%   50%   90%
## (Intercept)  0.3   0.8 -0.7   0.4   1.4
## acoustictness -0.4   0.1 -0.5  -0.4  -0.3
## danceability  0.3   0.1  0.2   0.3   0.4
## duration_ms   0.2   0.1  0.2   0.2   0.3
## energy        0.1   0.1  0.0   0.1   0.2
## instrumentalness 0.3   0.1  0.3   0.3   0.4
## liveness      0.1   0.1  0.0   0.1   0.1
## loudness      -0.4   0.1 -0.5  -0.4  -0.3
## mode1         -0.2   0.1 -0.3  -0.2   0.0
## speechiness    0.3   0.1  0.3   0.3   0.4
## tempo         0.1   0.1  0.0   0.1   0.2
## time_signature -0.1   0.2 -0.3  -0.1   0.2
## valence       0.2   0.1  0.1   0.2   0.3
```

```
##
## Fit Diagnostics:
##      mean    sd   10%   50%   90%
## mean_PPD 0.5    0.0  0.5   0.5   0.5
##
## The mean_ppd is the sample average posterior predictive distribution of the outcome variable (for de
##
## MCMC diagnostics
##      mcse Rhat n_eff
## (Intercept)    0.0  1.0  5838
## acousticness    0.0  1.0  5281
## danceability    0.0  1.0  4761
## duration_ms     0.0  1.0  6608
## energy          0.0  1.0  3811
## instrumentalness 0.0  1.0  6004
## liveness        0.0  1.0  6339
## loudness        0.0  1.0  4042
## mode1          0.0  1.0  6465
## speechiness     0.0  1.0  7079
## tempo          0.0  1.0  7195
## time_signature  0.0  1.0  5912
## valence         0.0  1.0  5250
## mean_PPD        0.0  1.0  4428
## log-posterior   0.1  1.0  1579
##
## For each parameter, mcse is Monte Carlo standard error, n_eff is a crude measure of effective sample
```




```
##      (Intercept)      acousticness      danceability      duration_ms
##          0.363          -0.413          0.311          0.228
##          energy instrumentalness      liveness      loudness
##          0.092          0.344          0.069          -0.422
##          model      speechiness      tempo      time_signature
##          -0.169          0.343          0.103          -0.056
##          valence
##          0.192
```

```
##              5%      95%
## (Intercept)   -0.997  1.687
## acousticness  -0.524 -0.305
## danceability   0.213  0.411
## duration_ms    0.136  0.323
## energy        -0.064  0.242
## instrumentalness 0.249  0.444
## liveness       -0.016  0.150
## loudness       -0.572 -0.274
## model         -0.331 -0.005
## speechiness    0.254  0.438
## tempo          0.019  0.188
## time_signature -0.391  0.292
## valence        0.096  0.294
```

```
#BIC
```

```
#View(spotify2)
```

```
#full_model <- stan_glm(target ~ acousticness+ danceability+ duration_ms+ energy+ instrumentalness+ liv
#              family = binomial(link = "logit"),
#              prior = cauchy(0,2.5), prior_intercept = cauchy(0,2.5),
#              seed = seed,
#              refresh = 0)
```

```
(loo3 <- loo(posterior3, save_psis = TRUE))
```

```
##
## Computed from 4000 by 2017 log-likelihood matrix
##
##      Estimate      SE
## elpd_loo -1267.9 16.3
## p_loo      14.4  0.5
## looic      2535.8 32.6
## -----
## Monte Carlo SE of elpd_loo is 0.1.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.
```

```
#Model Selection
```

```
##
```

```

## Computed from 4000 by 2017 log-likelihood matrix
##
##           Estimate  SE
## elpd_loo  -1398.9 0.5
## p_loo      1.0 0.0
## looic      2797.9 1.0
## -----
## Monte Carlo SE of elpd_loo is 0.0.
##
## All Pareto k estimates are good (k < 0.5).
## See help('pareto-k-diagnostic') for details.

##           elpd_diff se_diff
## posterior3    0.0    0.0
## posterior4 -131.1   16.3

## [1] 0.67

## [1] 0.661

```