

# The Top 100 Spotify Songs Each Year For 2010-2019

##Serena Everett (sle759) and Anika Pal (ap52553)

##Title and Introduction

```
## — Attaching packages ————— tidyverse 1.3.1 —
```

```
## ✓ ggplot2 3.3.5      ✓ purrr 0.3.4
## ✓ tibble 3.1.6       ✓ dplyr 1.0.8
## ✓ tidyr 1.1.4        ✓ stringr 1.4.0
## ✓ readr 2.1.2        ✓ forcats 0.5.1
```

```
## — Conflicts ————— tidyverse_conflicts() —
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
## Rows: 1000 Columns: 17
## — Column specification —————
## Delimiter: ","
## chr (5): title, artist, top genre, added, artist type
## dbl (12): year released, bpm, nrgy, dnce, dB, live, val, dur, acous, spch, p...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
## # A tibble: 6 × 17
##   title      artist `top genre` `year released` added    bpm  nrgy  dnce    dB  live
##   <chr>    <chr>   <chr>           <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 STARST... 3OH!3  dance pop           2009 2022...  140    81    61    -6    23
## 2 My Fir... 3OH!3  dance pop           2010 2022...  138    89    68    -4    36
## 3 I Need... Aloe ... pop soul           2010 2022...   95    48    84    -7     9
## 4 Airpla... B.o.B  atl hip hop         2010 2022...   93    87    66    -4     4
## 5 Nothin... B.o.B  atl hip hop         2010 2022...  104    85    69    -6     9
## 6 Magic ... B.o.B  atl hip hop         2010 2022...   82    93    55    -4    35
## # ... with 7 more variables: val <dbl>, dur <dbl>, acous <dbl>, spch <dbl>,
## #   pop <dbl>, top year <dbl>, artist type <chr>
```

The dataset we have chosen displays the top 100 songs across different genres for the years 2010 through 2019. We will use the year 2019 in this project because that year is very recent, so we know most of the songs. This dataset is interesting to us because we both use Spotify to listen to music and love using the app. This dataset was also found using the Kaggle website, where free datasets are available to the public. Source: (<https://www.kaggle.com/datasets/muhmores/spotify-top-100-songs-of-20152019?resource=download> (<https://www.kaggle.com/datasets/muhmores/spotify-top-100-songs-of-20152019?resource=download>))

*There are 1000 observations in the dataset and 17 variables in the dataset. The variables are: title of the song (title), artist of the song (artist), genre of the song (genre), the year the song was released (year released), the day the song was added to the Top Hits (added), beats per minute (bpm), how energetic the song is (nrgy), how easy it is to dance to the song (dnce), decibel (dB), how likely the song is a live recording (live), how positive the mood of the song is (val), duration (dur), how acoustic the song is (acous), the more the song is focused on spoken word (spch), popularity of song (pop), year the song was a top hit (top year), and if it is a solo artist or group (artist type). Because there are so many variables, we removed 7 variables (acous, live, spch, added, artist, val, and year released) we did not think were necessary for the project and kept the other 10.*

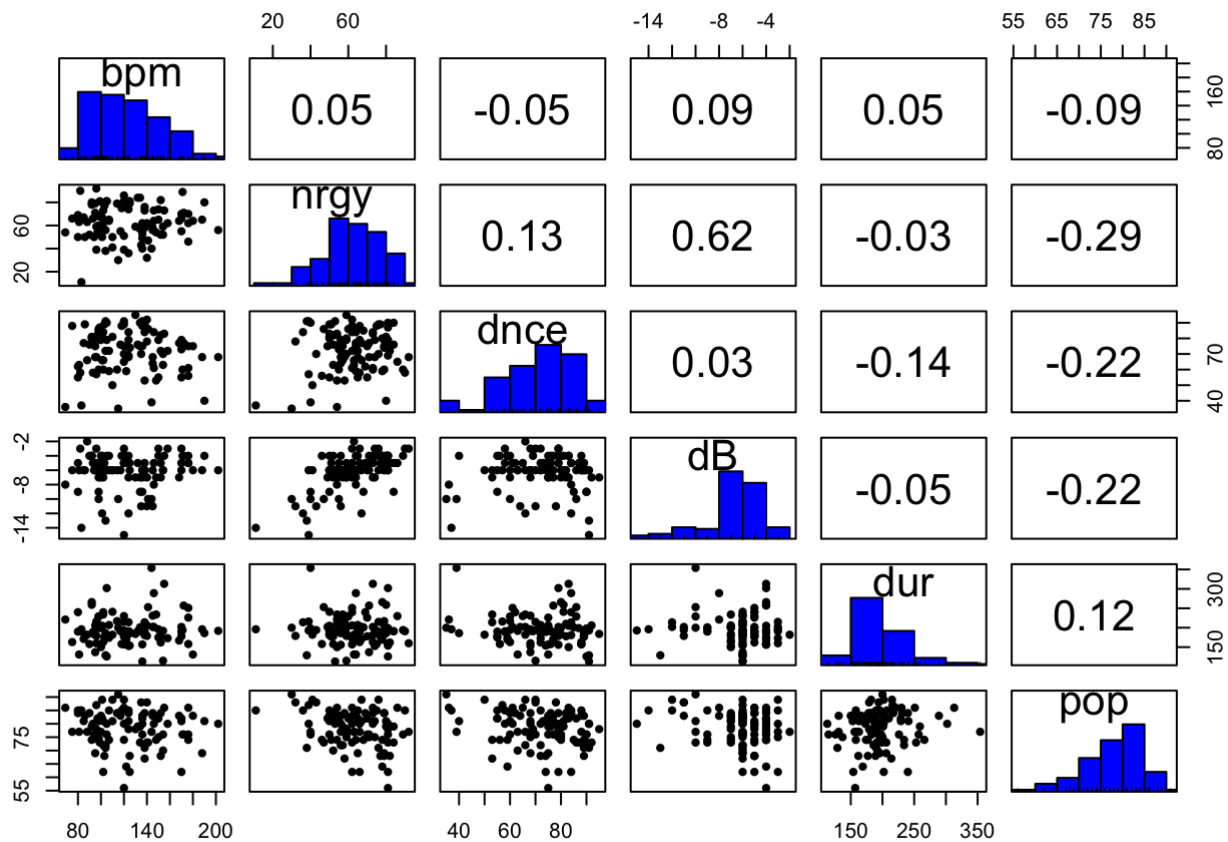
*Some possible relationships we expect are that the variables energy, danceability, and song popularity will have a positive relationship. This is because if the song is easier to dance to and is more energetic, the song may be more popular than others. There was no joining or tidying needed for this dataset as it was already tidy when we found it.*

*After importing the csv file into R, we renamed the dataset “spotify” so it would be easier to perform future functions. We saved the dataset with only 2019 songs as “spotify\_2019”.*

## Exploratory Data Analysis

```
##  
## Attaching package: 'psych'
```

```
## The following objects are masked from 'package:ggplot2':  
##  
##    %+%, alpha
```



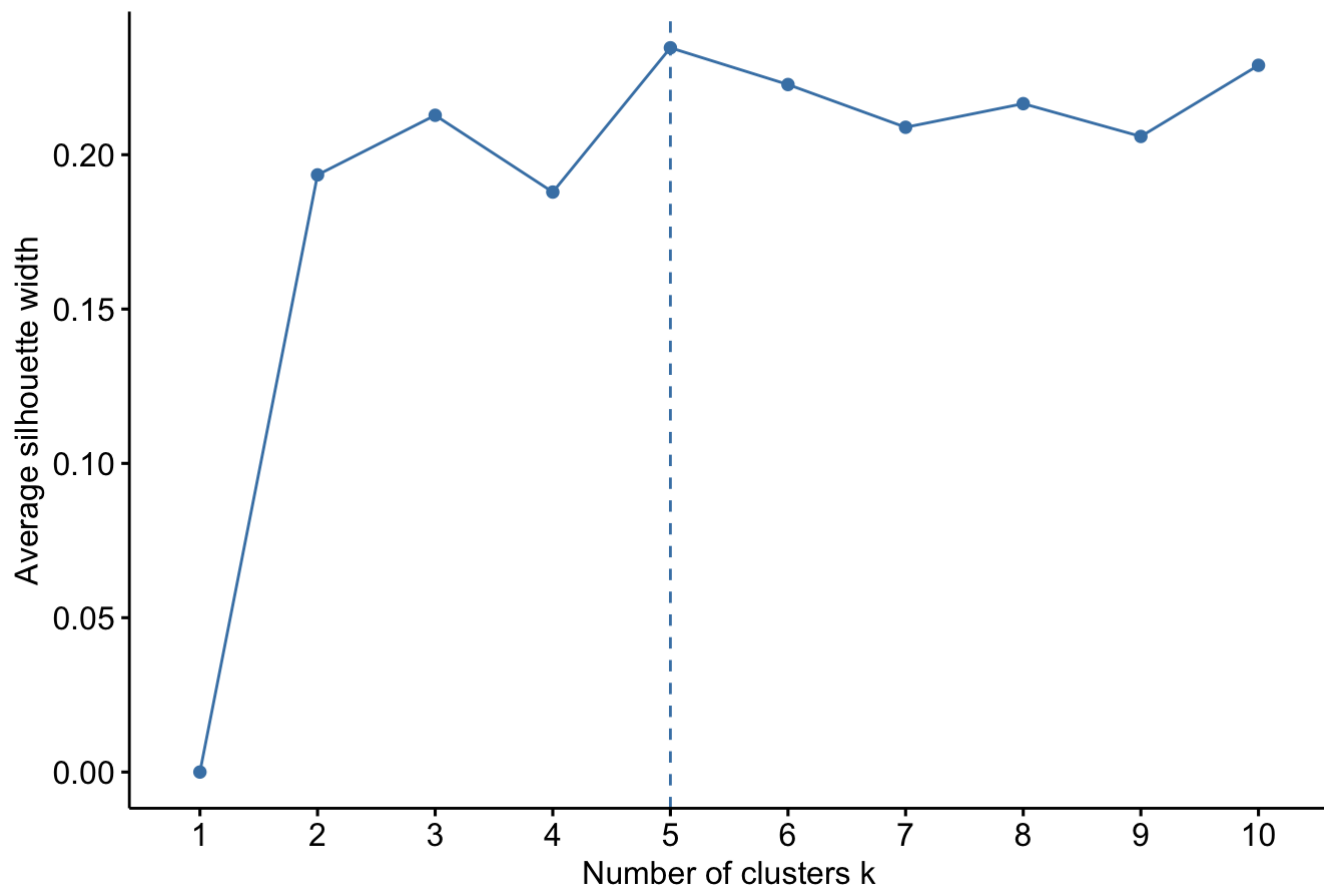
The

two variables that are the most correlated are *nrgy* and *dB* with a correlation coefficient of 0.62. The two variables that are least correlated are *dB* and *dnce* with a correlation coefficient of 0.03. There does not appear to be any relationships or trends that are apparent from the correlation matrix as there are no strong correlation coefficients. The *top.year* will appear as NA as a correlation and have a single point on the bottom graphs because the year is 2019 for all 100 songs. Therefore, *top.year* needs to be removed after we utilize it to filter for 2019 songs.

## Clustering

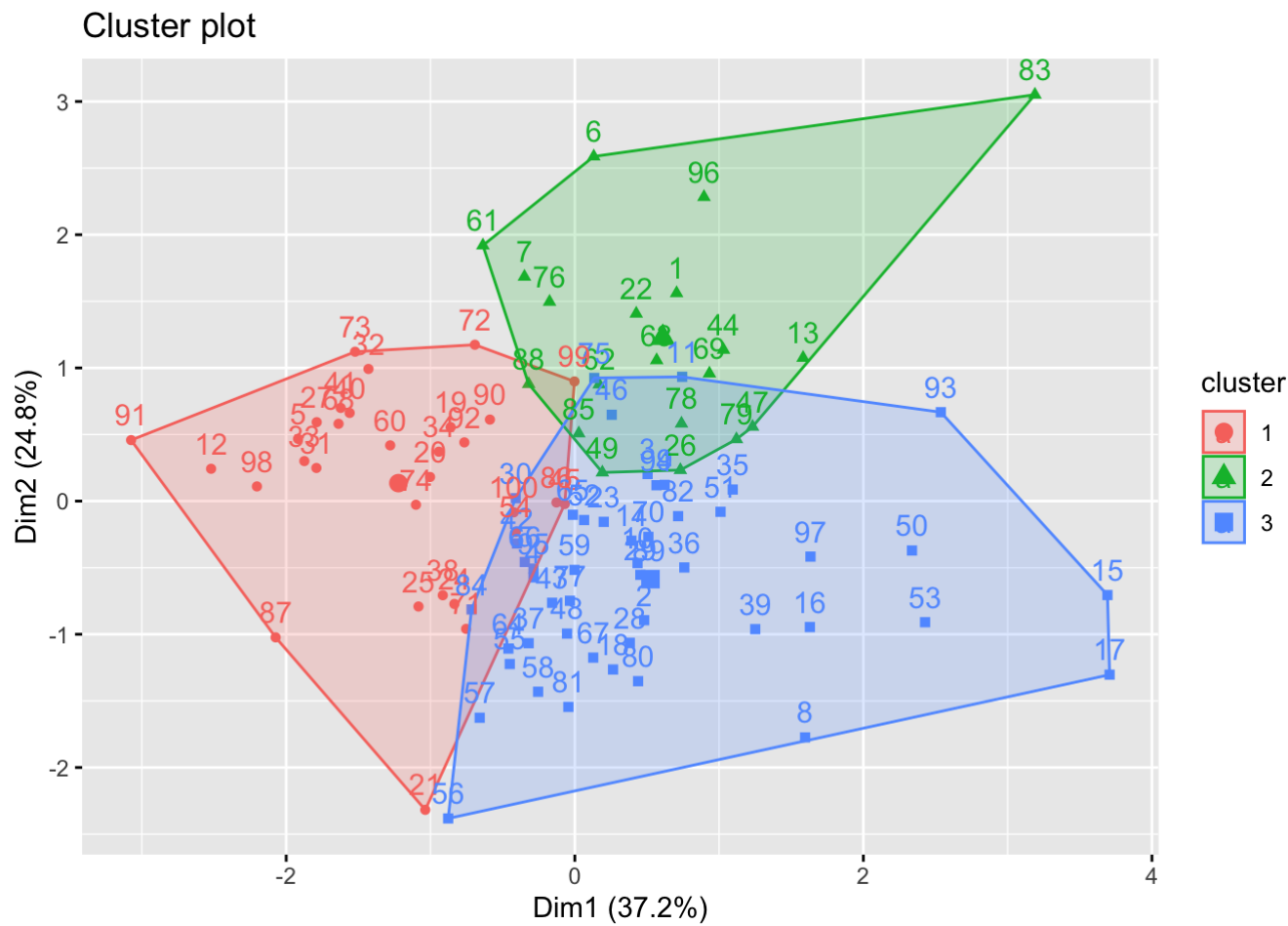
```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

## Optimal number of clusters



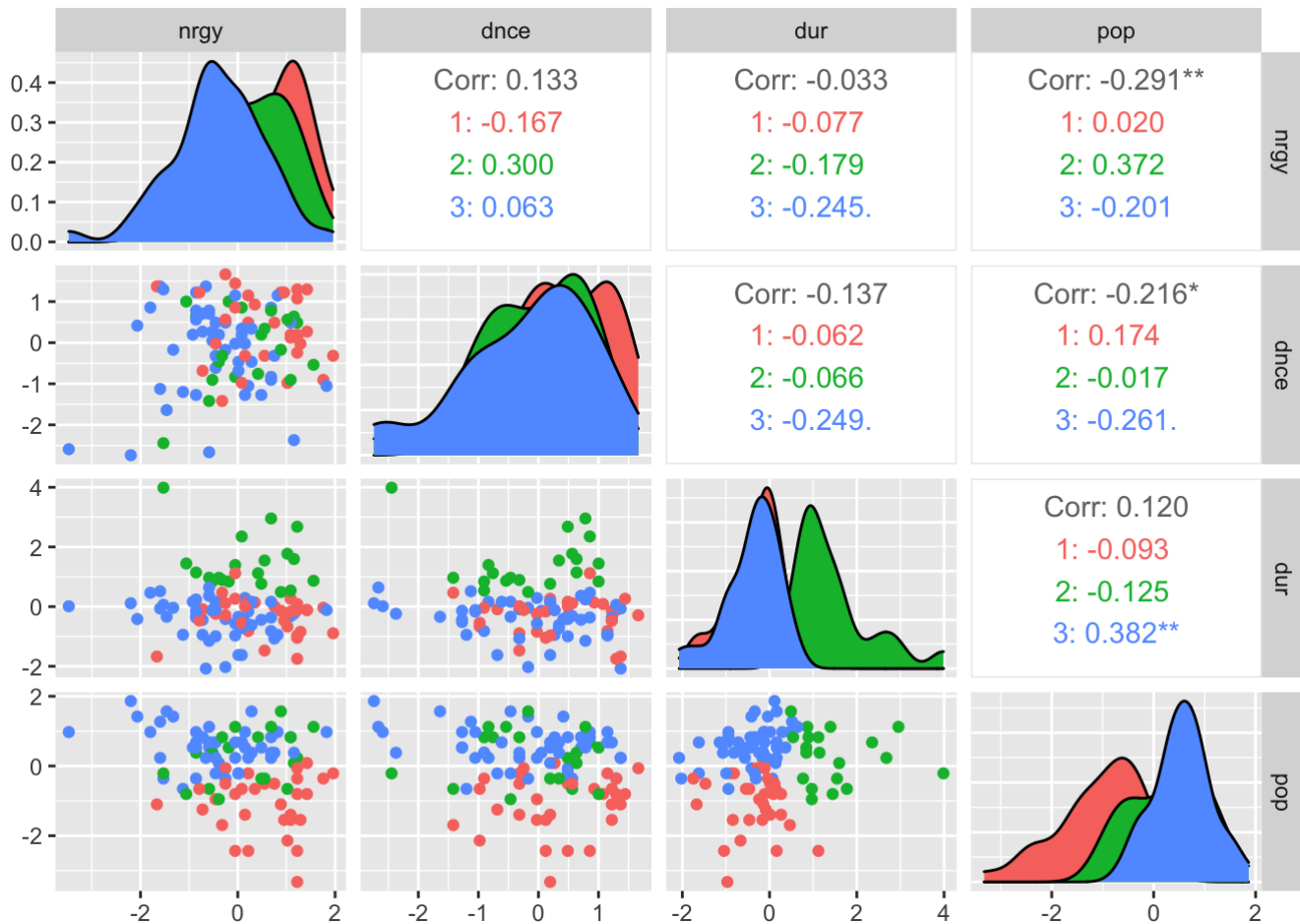
```
## K-means clustering with 3 clusters of sizes 31, 20, 49
##
## Cluster means:
##      nrgy      dnce      dur      pop
## 1  0.5399412  0.3536828 -0.3305408 -1.0771059
## 2  0.1817982 -0.1129722  1.4584353  0.2845131
## 3 -0.4157988 -0.1776474 -0.3861620  0.5653065
##
## Clustering vector:
##  [1] 2 3 3 3 1 2 2 3 3 3 3 1 2 3 3 3 3 1 1 1 2 3 1 1 2 1 3 3 3 1 1 1 1 3 3 3
## [38] 1 3 1 1 3 3 2 1 3 2 3 2 3 3 3 3 1 3 3 3 3 3 1 2 2 2 3 3 3 3 1 2 3 1 1 1 1
## [75] 3 2 3 2 2 3 3 3 2 3 2 1 1 2 3 1 1 1 3 3 3 2 3 1 1 1
##
## Within cluster sum of squares by cluster:
## [1] 75.48990 55.92731 134.25573
## (between_SS / total_SS = 32.9 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"
## [6] "betweenss"    "size"         "iter"         "ifault"
```

##	nrngy	dnce	dur	pop	cluster
## 1	0.08117187	0.8553606	2.3510982	0.67720051	2
## 2	-1.12634354	-1.1986786	-0.9429538	-0.65645476	3
## 3	0.21534025	-1.0519615	-0.1634453	0.23264875	3
## 4	-0.25424908	0.4885679	-0.4148997	-0.06371909	3
## 5	1.42285565	1.2955119	-0.1131545	-0.80463868	1
## 6	1.22160308	0.4885679	2.6779889	0.23264875	2



```
## # A tibble: 3 × 5
##   cluster  nrgy  dnce  dur  pop
##   <fct>    <dbl> <dbl> <dbl> <dbl>
## 1 1      0.540  0.354 -0.331 -1.08
## 2 2      0.182 -0.113  1.46  0.285
## 3 3     -0.416 -0.178 -0.386  0.565
```

```
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
```



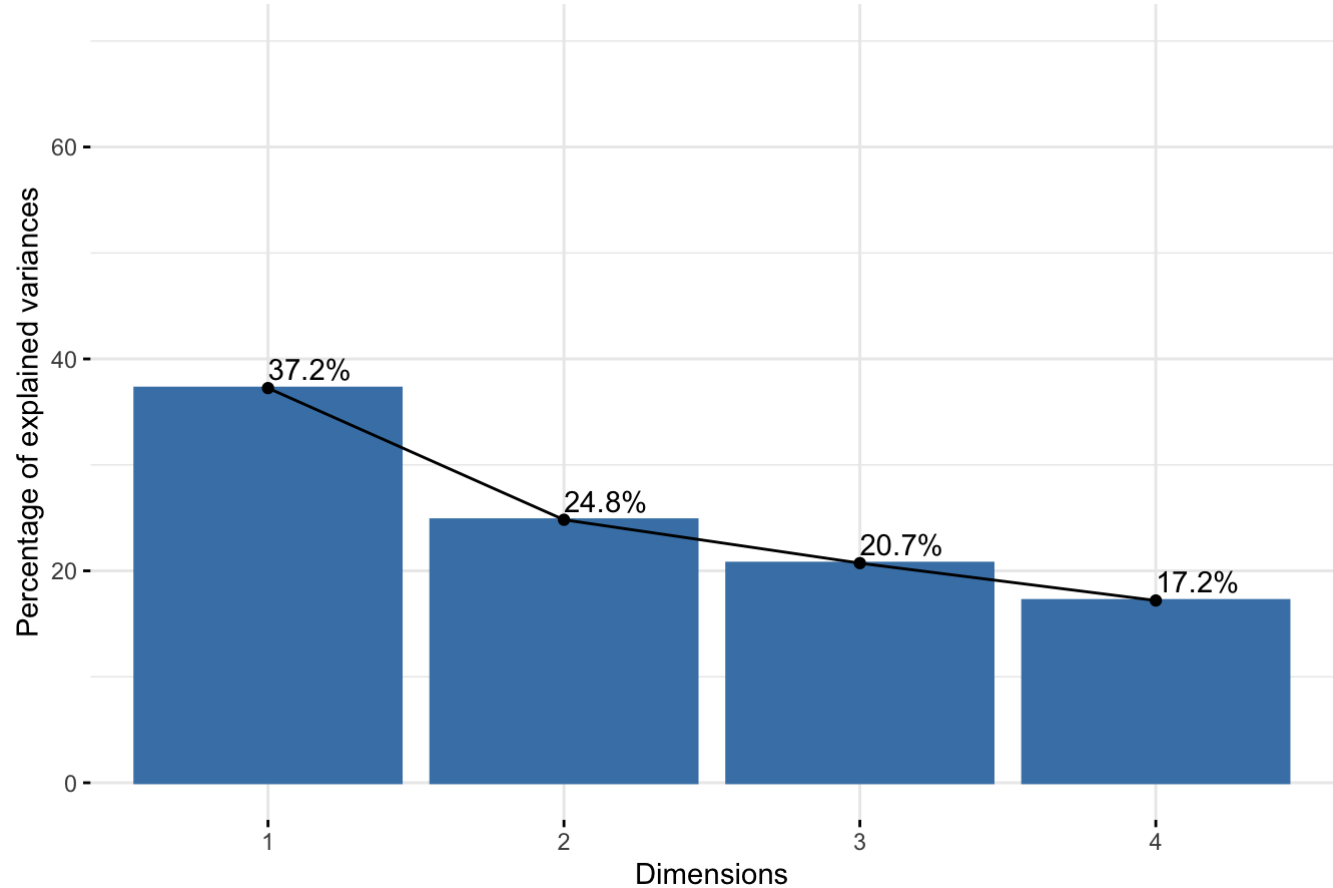
The `spotify_num` dataset, created in the previous section, needed to be scaled and was named `spotify_scaled`. After scaling the dataset, we found the optimal number of clusters using the silhouette width. Instead of using the optimal number of 5 clusters, Professor Guyot told us to use 3 clusters to help visualize the data better. The average silhouette width was around 0.25; this is not a high average, which indicates no substantial structure has been found. There are 3 clusters with sizes of 31, 20, and 49. We then mutated a column to help us identify which cluster each song was in. After creating a cluster plot, we see there is overlap in the center of the graph, which could mean that certain observations belong to more than one cluster. Therefore, the clustering is not very effective in separating the different songs. There might be some outliers in each group as well; for example, the green cluster extends so far out for the 83rd song. We then used a `ggpairs` plot to compare the correlation coefficients with the cluster distributions. None of the correlation coefficients were very strong between each of the variables. For the mean of each variable for cluster 1, `nrng` is 0.5399412, `dncc` is 0.3536828, `dur` is -0.3305408, and `pop` is -1.07769. The positive values indicate it is above the overall mean, and the negative values indicate that it is below the overall mean.

## Dimensionality Reduction

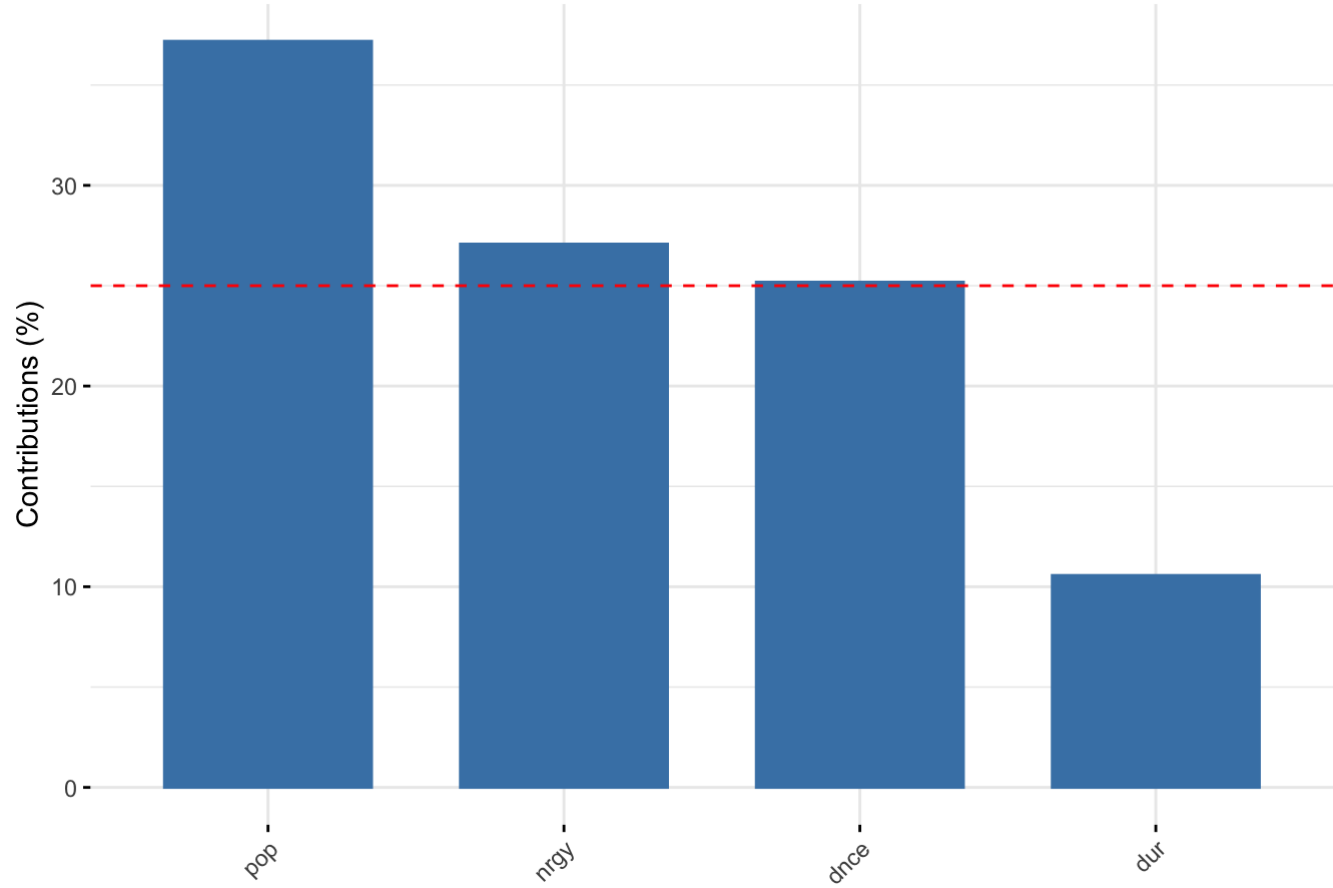
```
##      eigenvalue variance.percent cumulative.variance.percent
## Dim.1  1.4899495          37.24874          37.24874
## Dim.2  0.9928738          24.82185          62.07058
## Dim.3  0.8291173          20.72793          82.79851
## Dim.4  0.6880594          17.20149         100.00000
```



Scree plot

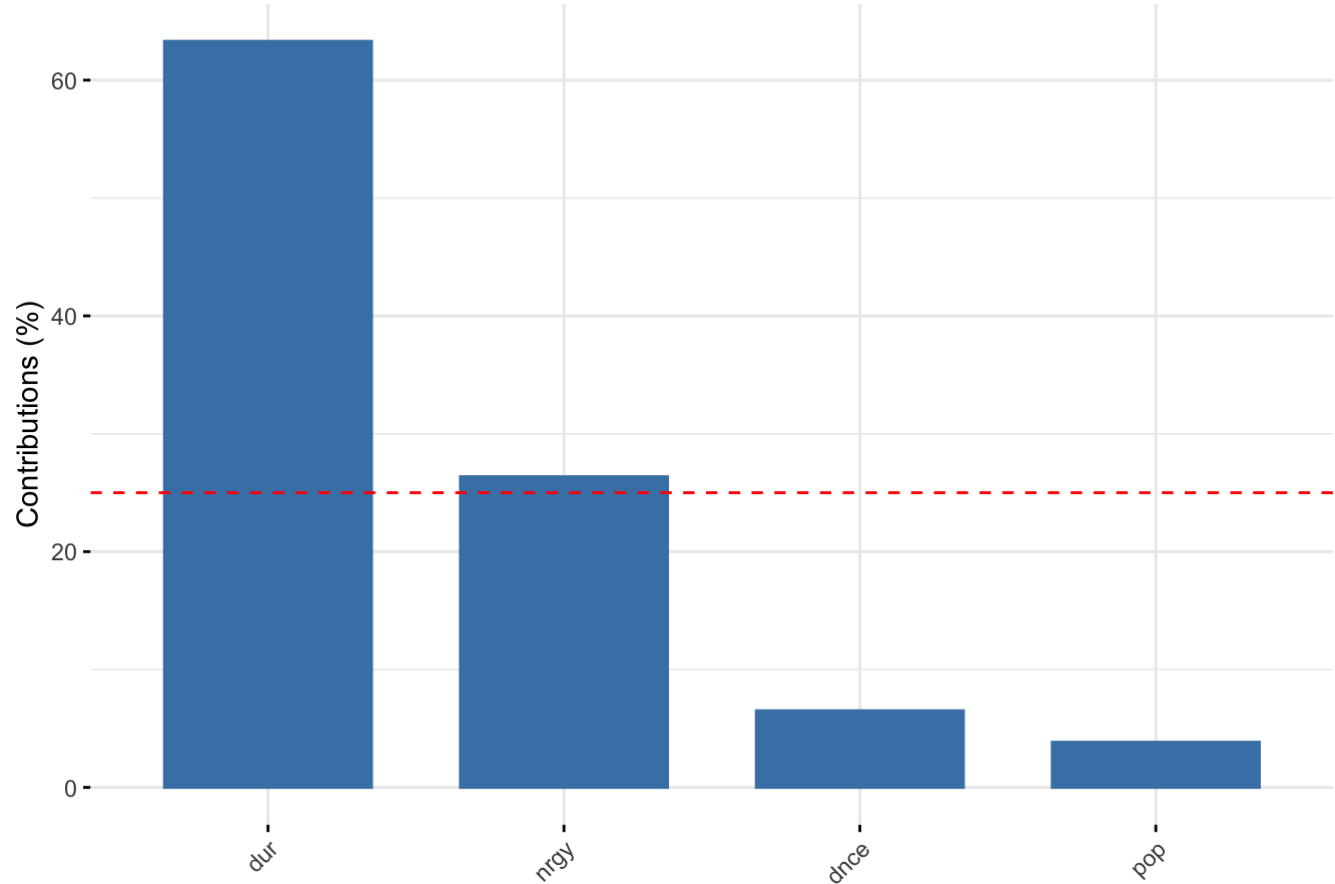


Contribution of variables to Dim-1

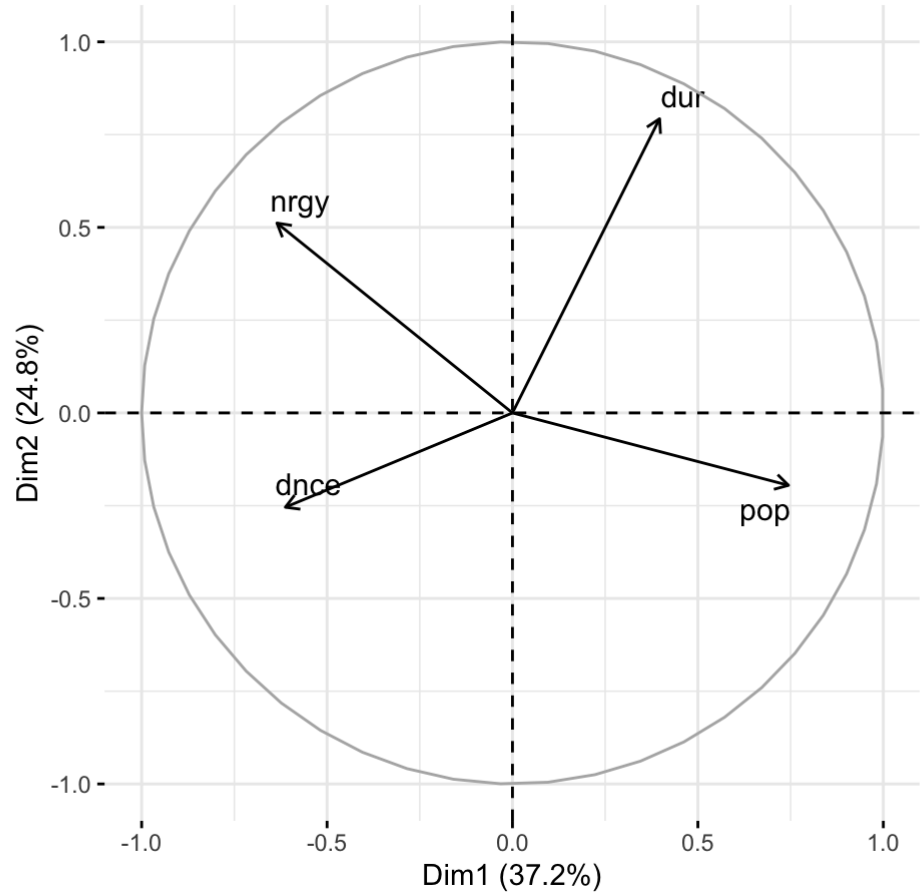


Contribution of variables to Dim-2





Variables - PCA



We

performed a PCA on four of our variables from the spotify\_scaled dataset, and we used a prcomp() function to find the principal components. The eigenvalue function gives us the number of principal components. When looking at

the PCs, we want to look at the first few that give us a total variance of roughly 80%. This means that we would use the first 3 PCs because the cumulative variance is 82%. This variance can be visualized by the scree plot. For the first PC compared to the second PC, there is less contribution difference between the variables NEED 2

ANSWER LAST 2 POINTS ON RUBRIC.

## Classification and Cross-Validation

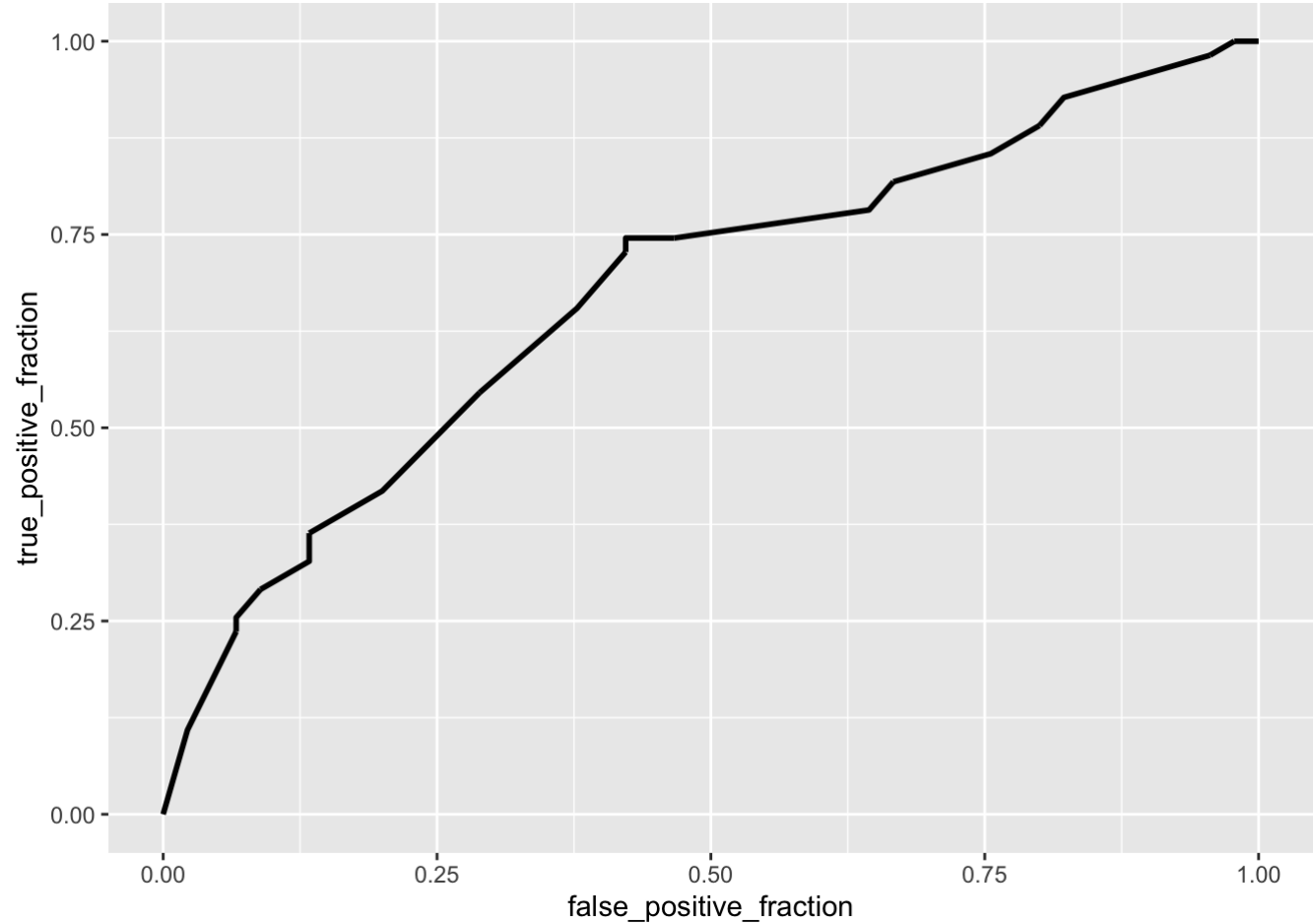
```
## Loading required package: lattice
```

```
##
## Attaching package: 'caret'
```

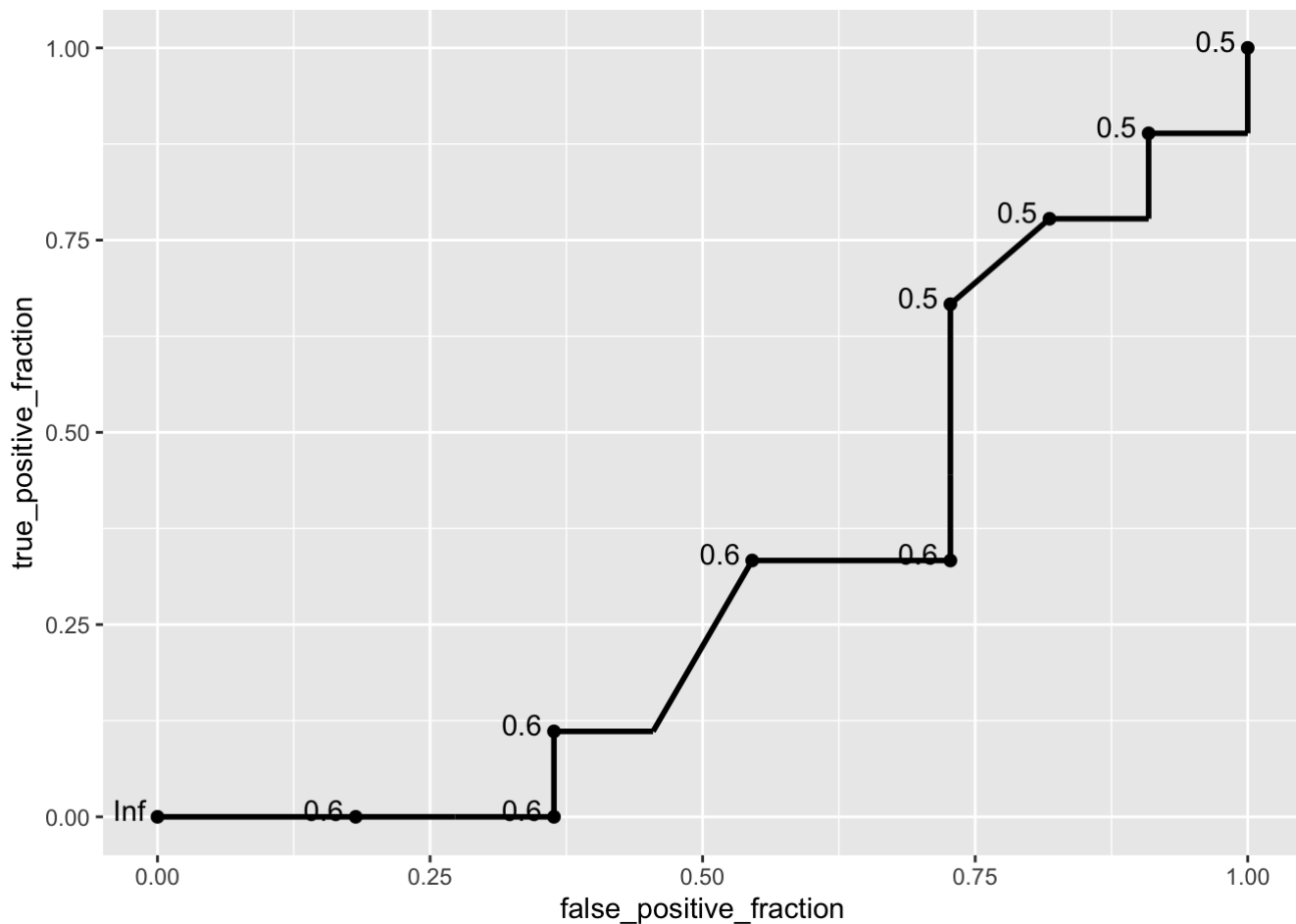
```
## The following object is masked from 'package:purrr':
##
## lift
```

```
## # A tibble: 1 × 1
##   `mean(dnce)`
##   <dbl>
## 1       72.3
```

```
## # A tibble: 6 × 5
##   title                nrgy danceability proportion predicted
##   <chr>                <dbl>         <dbl>         <dbl>         <dbl>
## 1 a lot                64             1         0.333             0
## 2 Easier               46             0         0.667             1
## 3 Swervin (feat. 6ix9ine) 66             0         0.375             0
## 4 Look Back at It       59             1         0.571             1
## 5 Ladbroke Grove       84             1         0.7               1
## 6 China                81             1         0.833             1
```



##	PANEL	group	AUC
## 1	1	-1	0.6731313



```
## [1] 0.4950712
```

Using the `spotify_2019` dataset, we found the mean of the danceability variable (`dnce`), and then we split the data based on the mean of 72. If an observation for `dnce` fell above the mean of 72, it was classified as 1. If it was below 72, it was classified as 0. Using `k` nearest neighbors functions, we said that if danceability is equal to 1 then it is a positive value and would show up as `TRUE`. Anything else would show up as `FALSE`. This was done with a `kNN` of 5. Then we created a proportion, called `kNN_spotify`, of the nearest neighbors with danceability and classified anything above 50% as 1 and anything below as 0. We then created a ROC curve based on the `kNN_spotify` proportion. The ROC curve was poor since it had a value of 0.673. We then performed a `k-fold` cross-validation with the same classifier and used 5 folds. This means the dataset was split into 5 parts. One part was considered the training set while the other part was the test set. This was repeated 5 times, so each part had been used once as a test set and the average performance was found after the 5 tests were performed. There is a noticeable decrease as the previous model had an AUC of 0.673 and the average five-fold AUC is 0.5020424. Therefore, the model does show signs of over fitting.