

Maddie Bauer  
DSC 550  
Original Case Study

**Topic:** Analyze Features of Songs from Spotify to Understand a Song's Popularity Score

**Hypothesis:**

What features of a song can lead to a higher popularity score on Spotify? Are any particular features of a song highly correlated to one another and/or have a significant effect on their popularity score?

**Background:**

Music plays a big part in many peoples lives, including mine. For this analysis, I wanted to explore data from my favorite music platform, Spotify, to understand what features of a song leads it to being more popular. Features of songs include items such as tempo, energy level, mood, loudness and more. Most people simply enjoy the music they are listening to without thinking in depth about the features mentioned above, which is why I wanted to explore this data to see if any of these features play an important role in determining the overall popularity of a song. This information can be useful for song producers, artists, music marketers and other professionals within the music industry.

**The Data:**

I found my data set on [Kaggle](#). It consists of 603 observations and 15 variables (listed below) for the most popular song titles between the years of 2010 and 2019.

*Song Number* – The song's number in a set

*Title* – Title of the song

*Artist* – Artist of the song

*Top Genre* – The genre of the track

*Year* – The song's year in the Billboard

*Bpm* – (Beats Per Minute) The tempo of the song

*Nrgy* – (Energy) The energy of a song – the higher the value, the more energetic song

*Dnce* – (Danceability) The higher the value, the easier it is to dance to this song

*dB* – (Loudness) The higher the value, the louder the song

*Live* – (Liveness) The higher the value, the more likely the song was recorded live

*Val* – (Valence) The higher the value, the more positive mood for the song

*Dur* – (Length) The duration of the song

*Acous* – (Acousticness) The higher the value, the more acoustic the song is

*Spch* – (Speechiness) The higher the value, the more spoken word the song contains

*Pop* – (Popularity) The higher the value, the more popular the song is

(Source: <http://organizeyourmusic.playlistmachinery.com/> )

## **Graph Analysis:**

**Step 1:** Load data into a dataframe

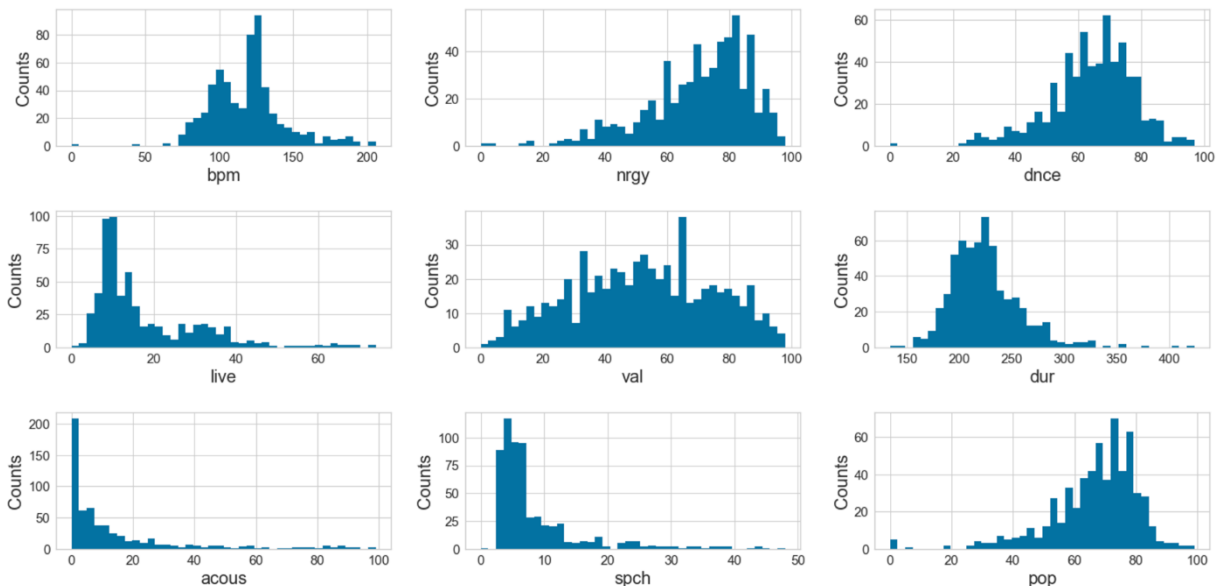
**Step 2:** Check the dimensions of the dataframe

**Step 3:** View the data

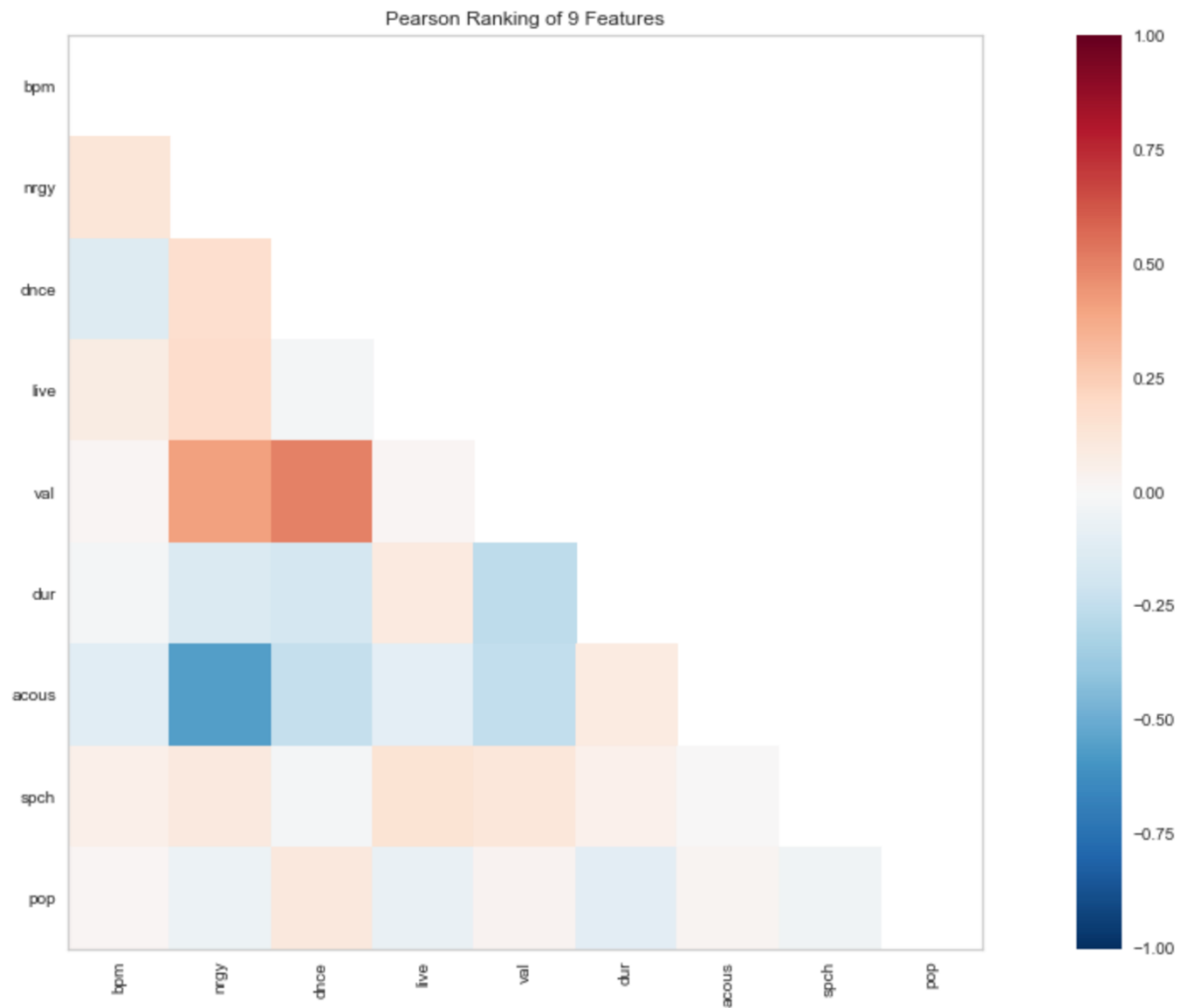
**Step 4:** Look at the different types of variables and summary information

**Step 5:** Clean up data (check for missing values, delete columns as needed, etc.)

**Step 6:** View distribution of variables with histograms



## Step 7: Pearson Correlation Ranking Chart

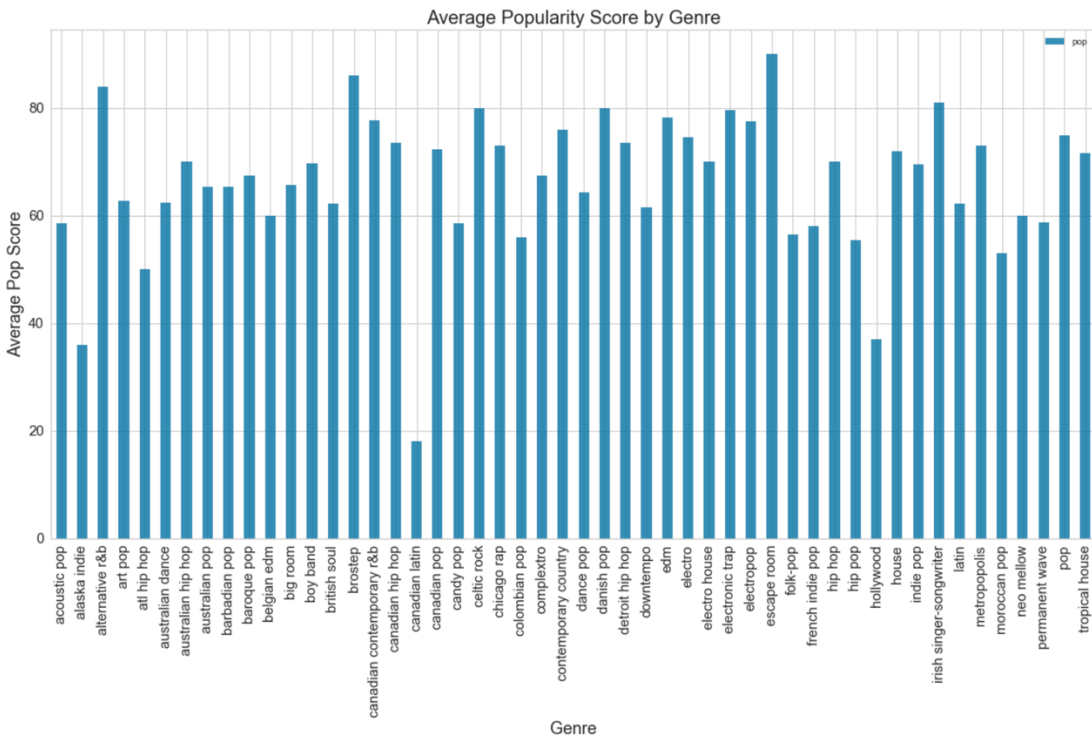


Dnce (danceability) and val (valence) are positively correlated.

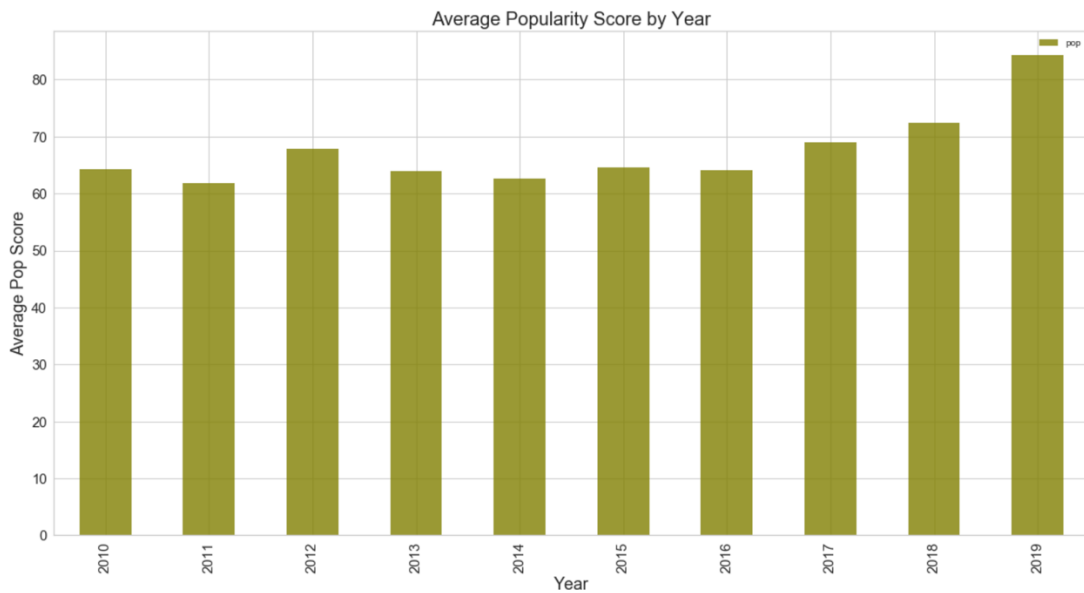
Acous (acousticness) and nrgy (energy) are negatively correlated.

There is no multicollinearity within this dataset.

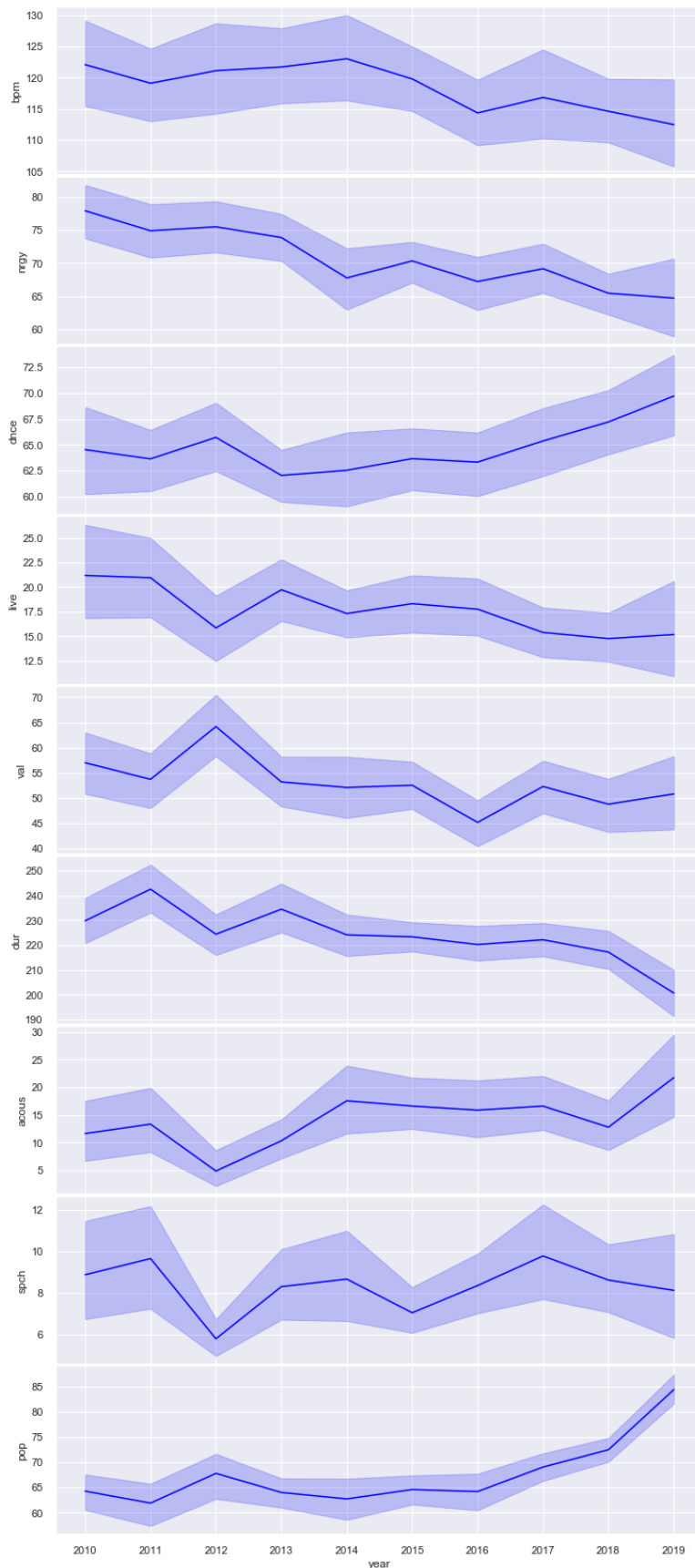
## Step 8: Explore some visualizations.



The genre with the highest popularity average is escape room. I have never heard of this genre before!



The year with the highest average popularity score average is 2019.

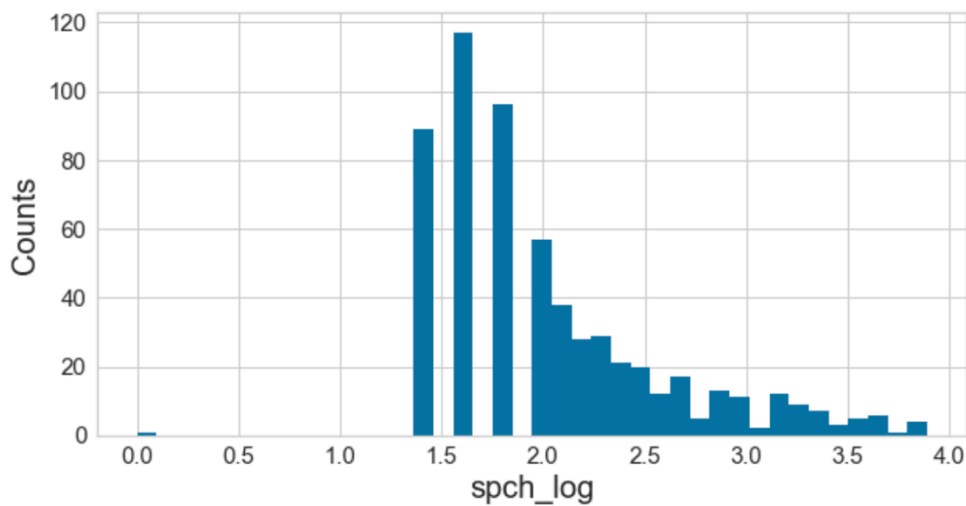
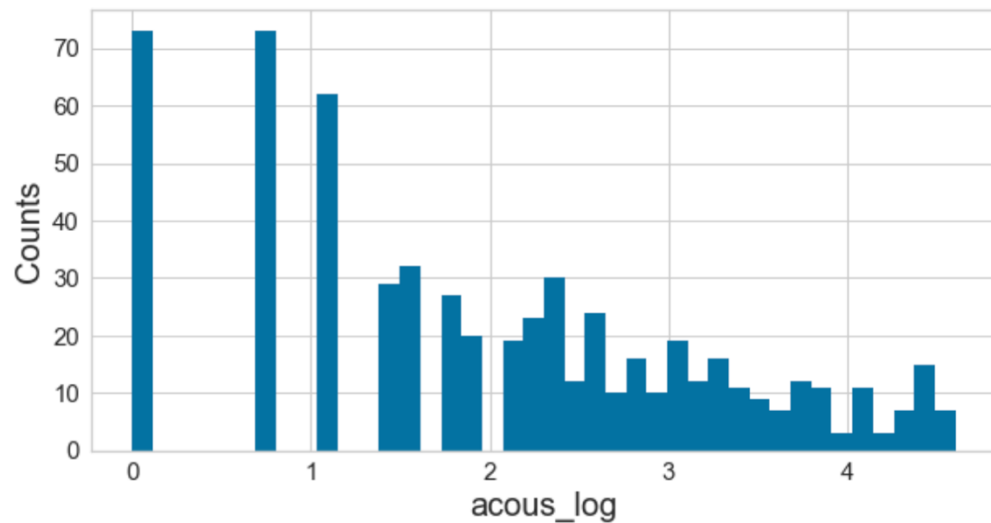


**From this plot I am noticing the following:**

1. bpm (tempo) is slightly decreasing over time
2. energy is consistently decreasing over time
3. danceability is consistently increasing over time
4. valence (mood) is consistent over time
5. duration is consistently decreasing
6. acousticness is consistently increasing over time
5. popularity is consistently increasing over time

## Dimensionality Reduction:

**Step 9:** Perform log transformation on skewed variables (acous & spch)



**Step 10:** View the variance of the features.

Several features have high variance which mean the spread in the data is quite large.

year	6.796747
bpm	614.809794
nrgy	266.037773
dnce	178.990088
dB	7.828917
live	171.676622
val	506.836091
dur	1164.860950
acous	431.233621
spch	55.997719
pop	210.764935
acous_log	1.609541
spch_log	0.335319

year	0.241261
bpm	0.018983
nrgy	-0.057645
dnce	0.116054
dB	0.156897
live	-0.075749
val	0.038953
dur	-0.104363
acous	0.026704
spch	-0.041490
pop	1.000000
acous_log	0.080561
spch_log	-0.005666

nrgy	-0.057645
dnce	0.116054
dB	0.156897
live	-0.075749
val	0.038953
dur	-0.104363
acous	0.026704
spch	-0.041490
pop	1.000000
acous_log	0.080561
spch_log	-0.005666

dB	0.156897
live	-0.075749
val	0.038953
dur	-0.104363
acous	0.026704
spch	-0.041490
pop	1.000000
acous_log	0.080561
spch_log	-0.005666

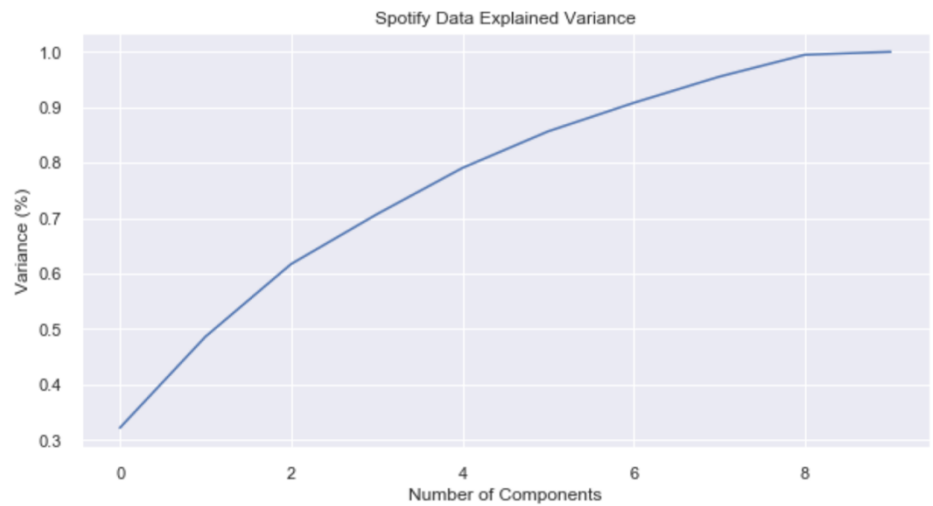
dur	-0.104363
acous	0.026704
spch	-0.041490
pop	1.000000
acous_log	0.080561
spch_log	-0.005666

```
acous_log    0.080561
spch_log     -0.005666
```

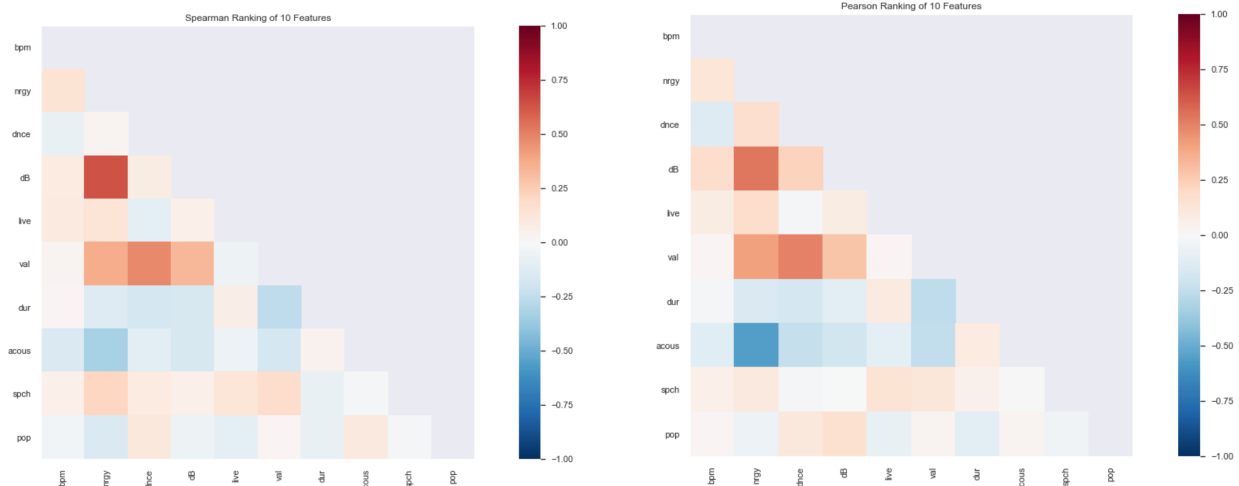
ar regression

[illegible]

**Step 14:** View the explained variance. The explained variance increases as the number of components (features) increases. I will use all 8 features.



**Step 15:** View Spearman and Pearson correlation heatmaps with normalized data frame.

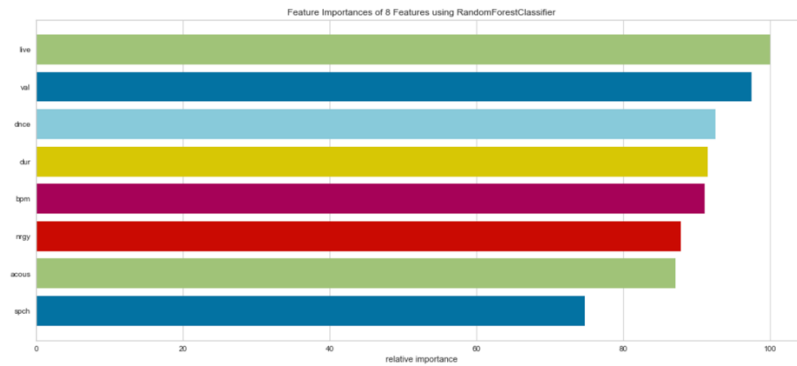


## **Model Evaluation & Selection**

**Step 16:** Split the data into two sets (training 80% and testing 20%).

**Step 17:** Create random forest classifier model and view feature importance





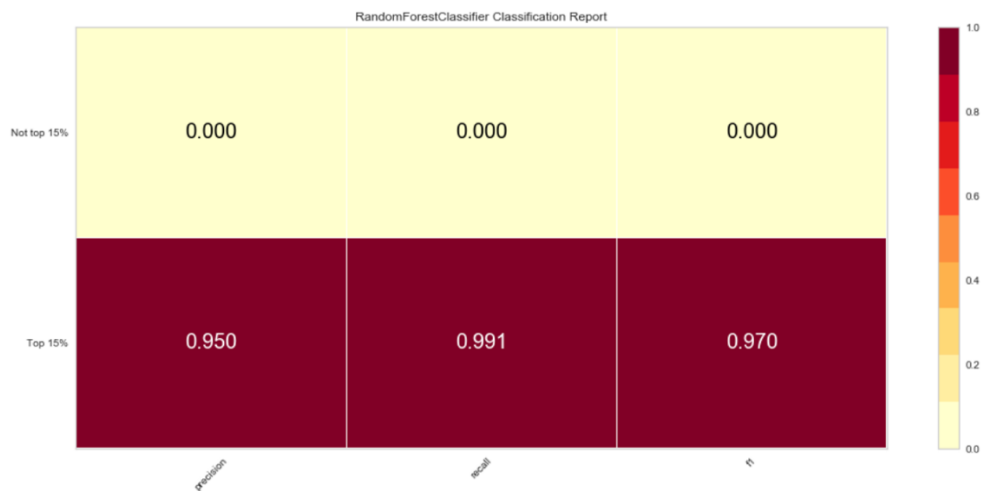
## Step 18: Display number of songs in the top 15% of each set and not in the top 15%

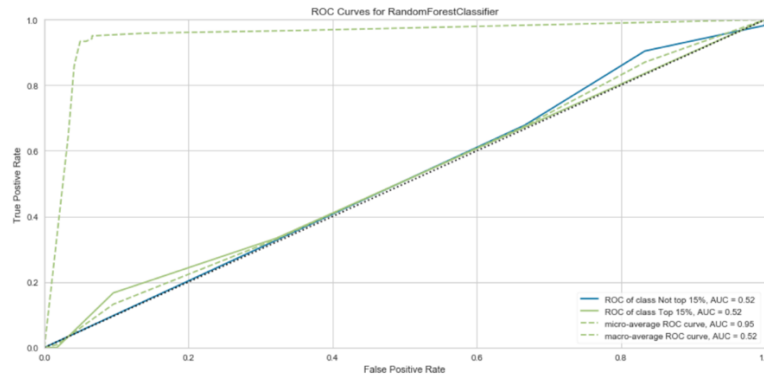
No. of samples in training set: 482  
 No. of samples in testing set: 121

No. of songs in the top 15% of pop in the training set:  
 Not top 15% 462  
 Top 15% 20  
 Name: pop, dtype: int64

No. of songs in the top 15% of pop in the testing set:  
 Not top 15% 115  
 Top 15% 6  
 Name: pop, dtype: int64

## Step 19: Classification Report and ROC





### Results of 1<sup>st</sup> Random Forest:

Precision: 0.950

Recall: 0.991

F1 Score: 0.970

**Step 20:** Try other models

### Results/Score

Support Vector Machine: 0.6363636364

Decision Tree: 0.925619834

Random Forest #2: 0.446280992

**Step 21:** Define and use function that will find similar songs to any randomly chosen song from the dataframe.

**Step 22:** Linear Regression

### Results of Linear Regression:

$$R^2 = 0.628$$

62.8% of the variance of the popularity score can be explained through the features of this dataset.

