

A Model for Predicting Pop Music Popularity and Its Different Characteristics Based on Multiple Linear Regression

Beinuo Guo

*United World College of Atlantic, Llantwit Major, Wales
Email: 2209886507@qq.com*

Abstract: *Pop songs are pretty diverse in the current digital music market. The article focuses on how different characteristics of a pop song can affect its popularity. In this paper, multiple linear regression is used to predict the model of pop song's popularity. Also, a Matlab code is made in order to achieve an ideal optimal popular pop song. The article can primarily answer the questions: What determines the popularity of a song? What kind of music do people like most currently? What characteristics shall composers focus on while making a new piece? This article may be helpful to those who make their own music and those who are engaged in the music market. Furthermore, This article also provides a computer model that can adjust parameters to obtain the optimal song type.*

Keywords: *model, predict, popularity, multiple linear regression, characteristic*

1. Introduction

Recently, pop music has occupied a large proportion of the music industry. Pop music is a genre of music that originated in the 1950s in the United States and the United Kingdom. Different from other genres with more sophisticated structures, such as classical music which includes forms like the sonata, pop music tends to be simpler and more repetitive. The rhythm of a pop song is also more diverse and novel than classical music.

In daily life, people often listen to pop songs. They listen to pop songs while they are drinking at a bar, listen to pop songs while driving a car, and listen to pop songs even while they are taking a shower. Plenty of young people are crazy about listening to all kinds of pop songs, and they are also fond of some special singers such as Taylor Swift, Ed Sheeran and so on. Different singers have different music tastes. Ed's music is often more inclined to folk style, while Taylor's is more energetic, talking about females rights more often. (Sloan 2021) These are all tags being stuck on them. Songs created by the same singer are also quite different from each other.

Obviously, we can feel that the popularity of different songs is quite different from each other. Some are praised by people all over the world while a few people only know some. Take Ed Sheeran as an example. His pop song "Shape of You" is quite popular worldwide, while he also wrote some other songs are not that popular. Why? Since the singer is the same, the main difference must be the music itself.

It is pretty hard to describe a piece of music in text, but a few characteristics can indeed summarize it. Some characteristics, such as the song's length and the published date of the song, are quite apparent, while others, such as speed and contrast, are entirely subjective. Researching both characteristics of a song helps give a more precise answer to the question above.

In this article, in order to eliminate the difference of region, all the songs collected are in English. Also, songs' polarities might also be affected by its record company and its propaganda. However, due to the data restriction, we only consider the song itself and the singer. As a result, this article might have some small errors comparing to the reality. However, it does help the music producer discover the public's tastes.

This article collects diverse data of 150 songs published on YouTube, including publication date, the number of followers of the YouTuber, and length of the song. Also, the article estimates the index of speed, contrast, and beat of each song. After processing the raw data, a model of how the characteristics of a pop song related to its popularity is shown in the article. Section II mainly talks about some

additional information in other articles about pop songs; Section III shows the data is collected and processed in order to make the model more precise; Section IV gives the model and estimates the effect of the model, with some statistical analyses on the data, while Section VI offers some conclusive comments and results.

2. Literature review

Previous work has given some analysis on music popularity. Many works have shown that some characteristics can affect music popularity. Some articles used a computer algorithm to show how different songs are spread. The following figure shows how different characteristics of a song relate to each other. (Bernardo, 2017)

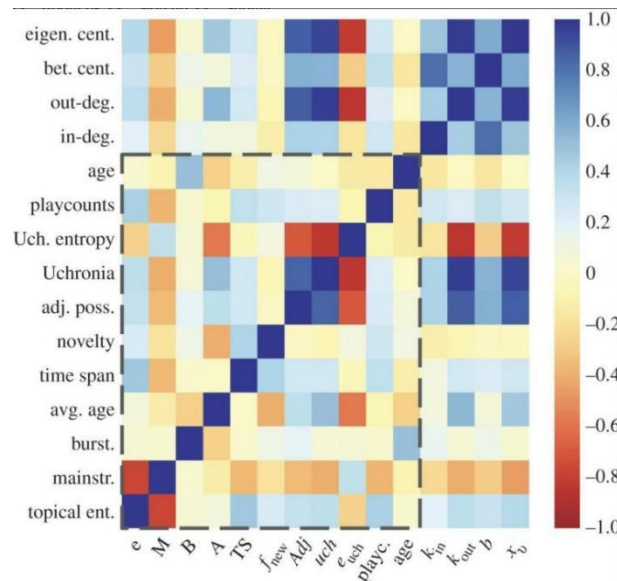


Figure 1: Coefficient Correlation of different characteristics matrices of pop song. Significance and popularity in music production, Bernardo

The color band on the right side represents the degree of coefficient correlations between two of the indices collected in the article: the strong anti-relation of Uchronia Entropy (Uch. Entropy) and Uchronia following the anti-relation of the mainstream index (mainstr.) and topical entropy index (topical ent.). Since they are related to the long-term effects of mainstream (mainstr.) and topical entropy (topical ent.), these two strong anti-relations illustrate that small topical diversities are usually related to specialized 'niche' albums, far from mainstream production.

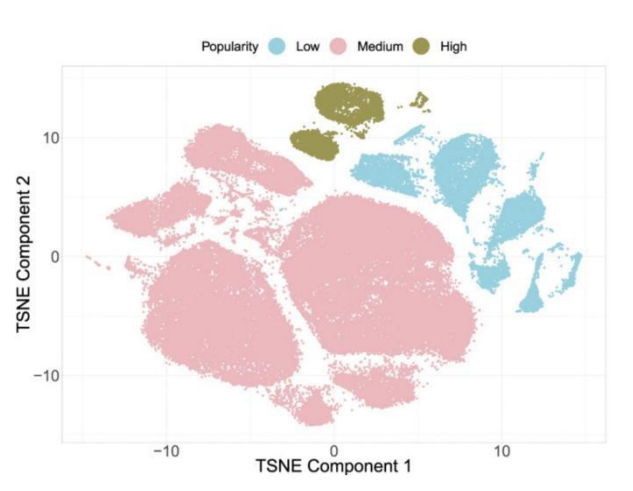


Figure 2: TSNE representation of the compressed feature x^T regarding the popularity level: Low (Orange), Medium (Purple) and High (Green). A Multimodal End-to-End Deep Learning Architecture for Music Popularity Prediction. David Martin-Gutierrez

Other articles transfer the songs into a different feature vector and calculate plenty of indexes such as Acousticness, Danceability, Duration, and so on in order to get a more precise result. (David Mart ́n-Guti ́rez, 2020)

Previous work also analyzes how a single variable can affect popularity. We can estimate that liveness does have negative effects on popularity. However, since the lower bar of speechiness is specially dense and focused, it is hard to estimate the relationship between popularity and speechiness. After analyzing plenty of indexes, the conclusion is that Speechiness, Instrumentalness, and Live are the features that negatively affect the Popularity Index, while Energy, Valence, and Duration of the song are the ones that positively affect it. (Sciandra 2020, Mart ́n-Guti ́rez, 2020)

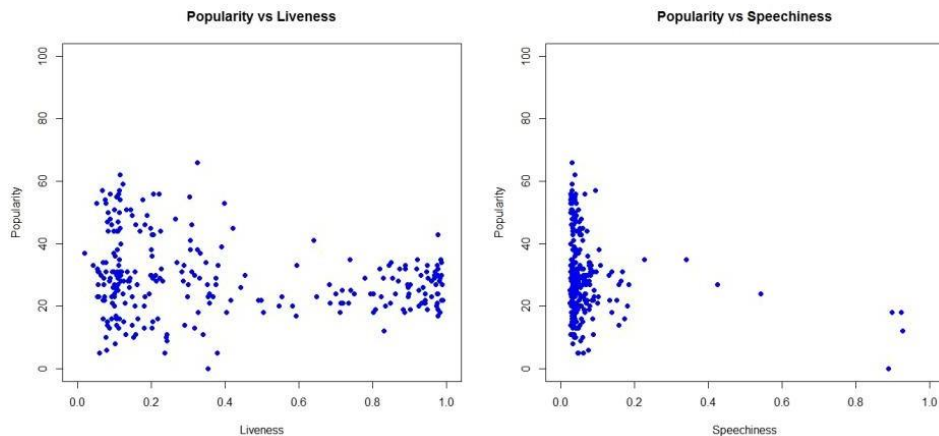


Figure 3: Scatterplot Popularity vs. Liveness and Speechiness A model-based approach to Spotify data analysis: a Beta GLMM, Sciandra

3. Data collection

The data set is collected from YouTube, a video application. While opening a song, some basic information such as published date, followers, numbers of viewing, and so on appears on the page. This article records data mainly from the data that appears on YouTube.

Furthermore, I listened to all of them and gave them an estimated score on Beat, Contrast, and Speed to describe their characteristics. Below is my estimated model of music popularity.

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 + \epsilon$$

Y_i is the index of popularity, calculated by total viewing times divided by published time. This index is the average viewings per month. Since this article aims to find out the relation between a song's popularity and its characteristics, only looking at its total viewing times is not proper. There is a nearly linear relation between published date and viewing times since people always have access to listen to the song on You-Tube after the song was published.

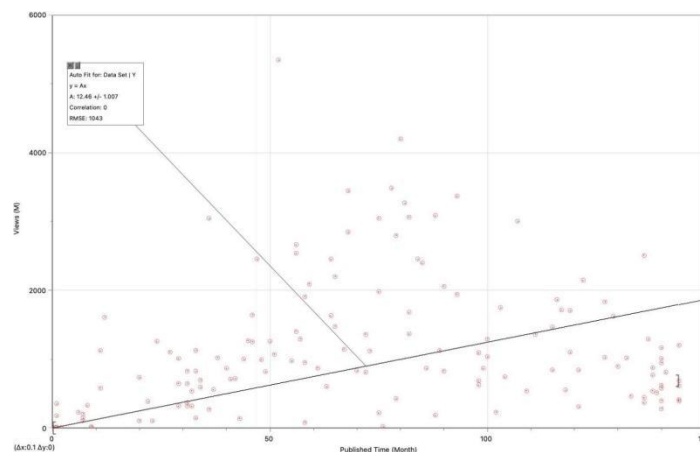


Figure 4: Relation between Published Time and Views

There are 4 total data recorded from the page: Length of the Song, Published time of the song, Follower of the You-tuber, and the numbers of viewing.

Length of the song represents how long the song is. The unit here is in seconds.

The published date is how long the song has been published on YouTube. In the article, this index is calculated by the difference between the current date and published date, with the month's unit.

Follower is the number of followers who subscribe to the YouTuber. This index can measure the popularity of the singer, which somehow effects the popularity of the song. Additionally, this article also estimates some internal characteristics of songs. Three additional indexes are being estimated and collected in the data collection, which are speed, contrast, and beat. They are all estimated and graded on a scale of 1-10.

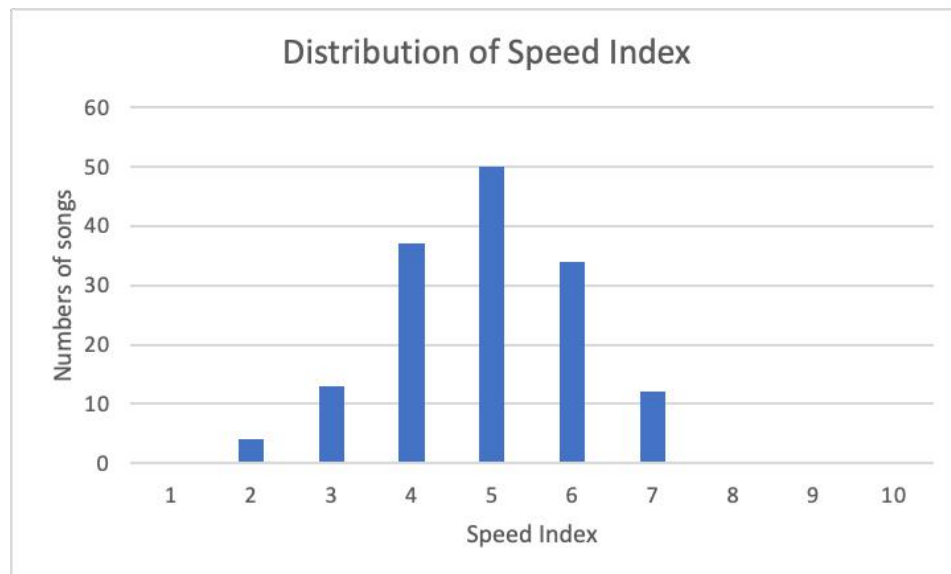


Figure 5: The Distribution of Speed Index

Figure 5 shows the distribution of the speed Index. Speed mainly measures how fast a song goes. Usually, it is measured by how many quarter notes can be played in a minute. Since it is hard to measure the precise speed of each pop song (pop songs usually have a quite unstable speed), the article estimated it using the method of grading. Take the song Shape of You, which gets a 7 for its grade. The song is at a fast speed from beginning to end. As a result, we include this song in the highest class, which is Grade 7.

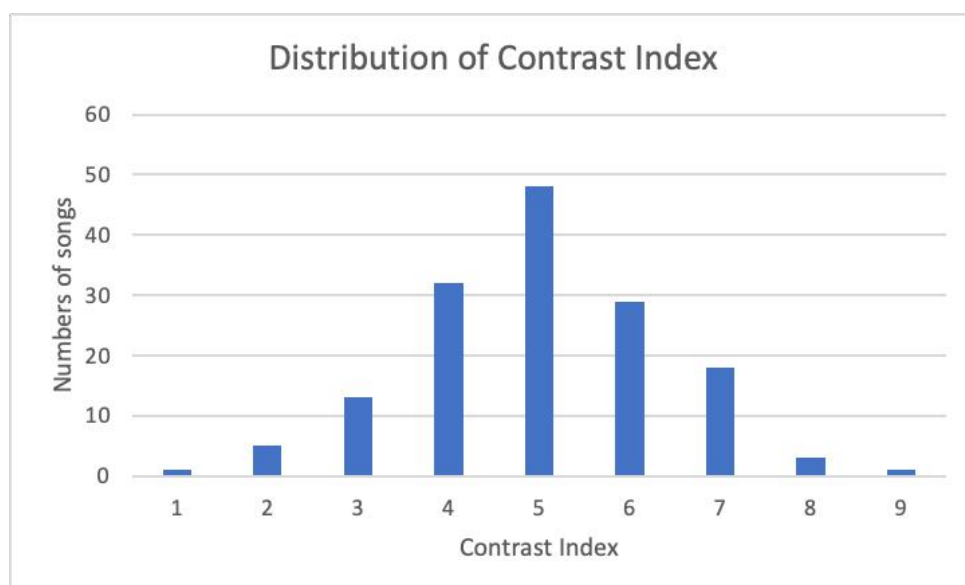


Figure 6: The Distribution of Contrast Index

Figure 6 shows the distribution of the contrast Index. In the diagram, there are 48 songs being graded 5 out of 10, which is also the most among all the grading.

The song Shallow by Lady Gaga and Bradley Cooper scores 9 out of 10 in this index. The song is sung by a man first, quite quiet and peaceful at the beginning, and the vocal range is in the same octave. Then, it changes to Lady Gaga using a female voice to sing the song but in a similar melody. Later, the vocal range suddenly rises dramatically. The beat also becomes more vital in the late part of the song. As a result, although the beat is not quite strong through the song, such a contrast of singers, vocal range, and beat give this song a grade of 9 out of 10.

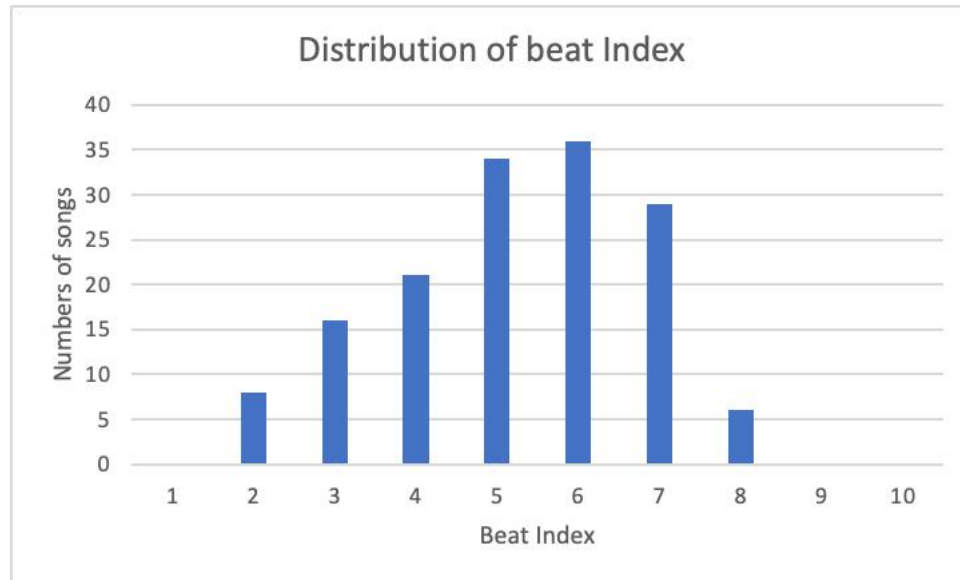


Figure 7: The Distribution of beat Index

Figure7 measures the sense of the rhythm of a song. It is quite different from the speed Index since the speed Index measures how fast the rhythm goes. A strong sense of rhythm can be like a strong base note and beat sound with a less elegant and continuous melody. The song Everybody, for example, scored 8 out of 10 on this index since there is a quite strong bass beat. Also, the accent of each sentence makes the song much more energetic and full of a sense of rhythm.

In probability theory, a normal distribution is a type of continuous probability distribution for a real-valued random variable. This distribution is quite proper for grading since the curve is like a bell, with few extreme grades and plenty of normal grades. The general form of its probability density function is:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Where $f(x)$ is the function of probability variables distributes, σ stands for the standard deviation, and μ stands for the expectation of the distribution of so-called mean.

Looking at the Speed data, we can make a statistical analysis and find that the total number of observations is 150 with an average Speed Index of 4.887 and a standard deviation of 1.17. Then we can substitute the value into the formula and get the $f(x)$ of the Speed Index.

$$f(x) = \frac{150}{\sqrt{2\pi} * 1.17} \exp\left(-\frac{(x - 4.887)^2}{2 * 1.17^2}\right)$$

The curve of the function shows as Figure 5:

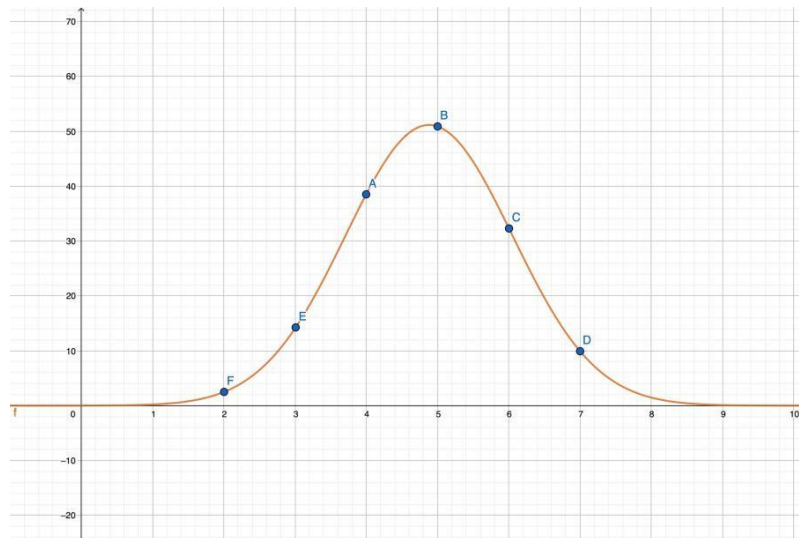


Figure 8: The Normal Distribution of Speed Index

The x-axis shows the Speed Index grading.

Point A represents the number of songs that get a four as the Speed Index assumed by the normal distribution. Point A on the graph has a y-value equal to 38, while the actual number in the data set is 37. Point B on the graph has a y-value equal to 51, while the actual number in the data set is 50. Basically, they are quite similar, and the data Speed Index mainly follows the normal distribution.

Implement the same step for Contrast and Beat Index:

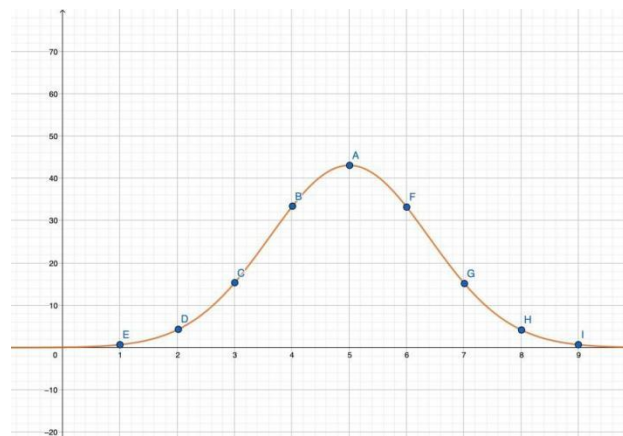


Figure 9: The Normal Distribution of Contrast Index

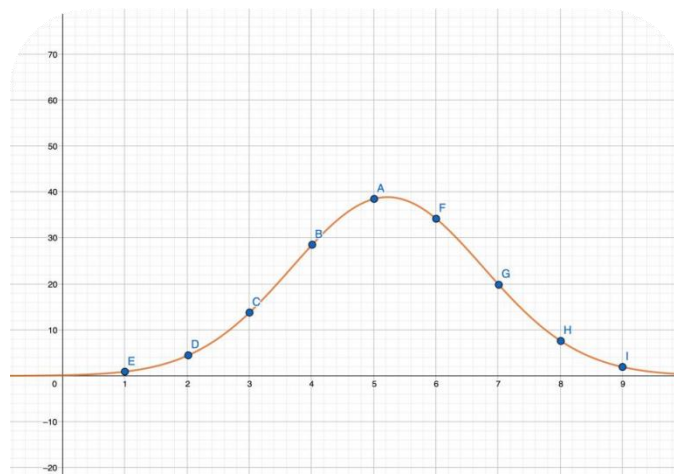


Figure 10: The Normal Distribution of Beat Index

Basically they follow Normal distribution which a quite scientific grading method.

3.1. Data process

Since this article uses multiple linear regression, one variable can only lead to the popularity going up or down when the variable is increasing. However, the large Speed Index and Beat Index mean the song is quite energetic and is suitable for dancing or cheering, while low Speed Index and Beat Index mean the song is quite elegant and suitable for chilling.

We cannot simply infer that people have a preference for the two types of pop song. As a result, this article does some data processing in order to show the type of music.

3.2. Intensity and Elegance

In order to measure the type of pop song, the article introduces two new variables called Intensity(I) and Elegance(E). The larger each variable is, the more intense (or elegant) the song is. Basically the more intense the song is, the less elegant the song is.

$$E = (10 - b) * (10 - c) * (10 - s)$$

$$I = b * c * s$$

3.3. Extremity

Since Intensity and Elegance still cannot solve the problem of linearity, the article introduce one more variable called Extremity(Ex). If the song is attributed to the area of high intensity or high elegance, extremity's value shall be high; on the contrary, it shall be low when the song is in the area of poor intensity or poor elegance.

As a result, we can use tanh function to fix the problem.

$$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} + 1$$

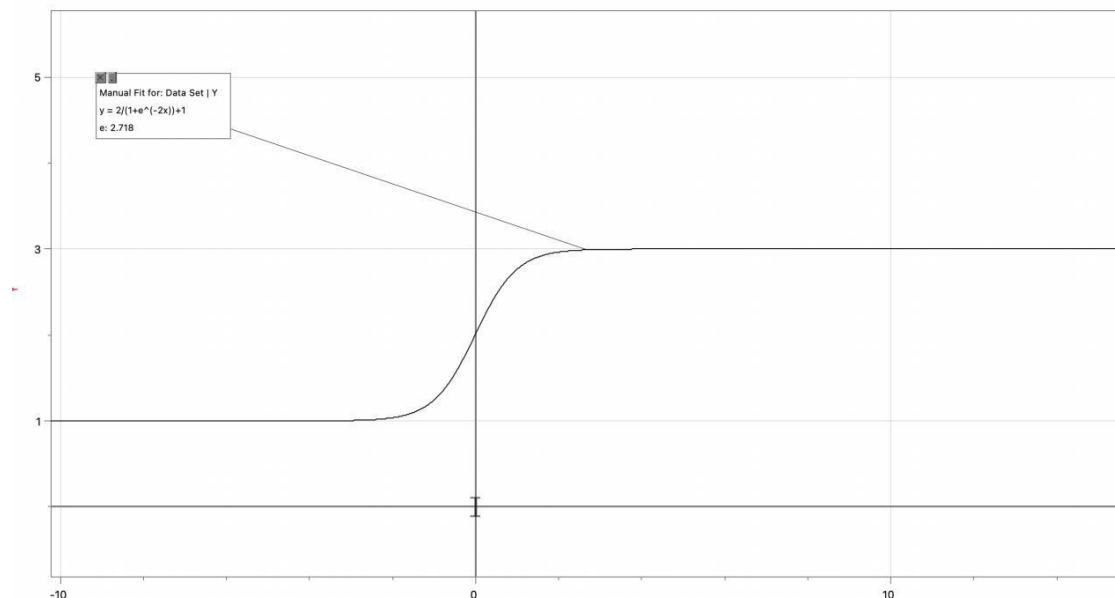


Figure 11: The TANH Activation Function Curve

Here x shall equal the absolute value of Intensity minus Elegance. The larger x is, the more extreme the pop song is. Since the average value of the extremity is 126, we aim to let the inflection point at x=126. Also, the value of Extremism shall be positive. What's more, the curve shall not be quite steep since we aim to give each different song a different Extremity. Thus, we made such variations:

$$Ex = \tanh(0.02a - 2.52) + 1 = \frac{2}{1 + e^{-0.04a + 5.04}}$$

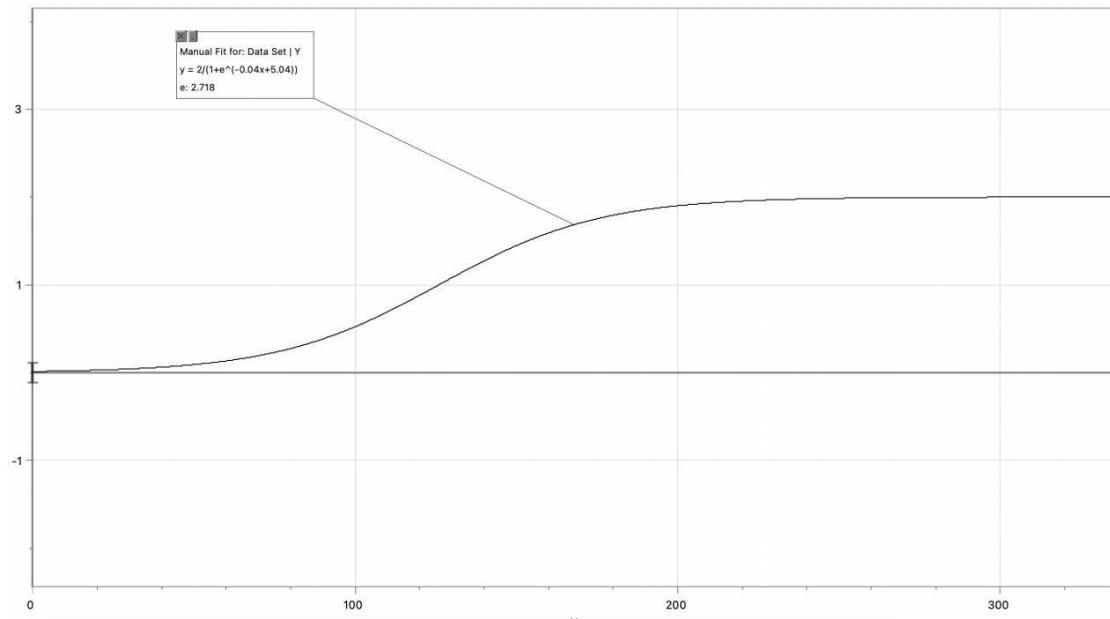


Figure 12: The TANH Activation Function for Extremity

We can clearly see that the songs that have a strong difference between intensity and elegance larger than 200 (strong extremity) basically share the similar extremity categorization since they are all labeled as “strong extremity” and the only difference might be some tiny difference in grading. (Same for the lower part)

Table 1 Descriptive statistical analysis of each index

| Index | Average | Median | Standard deviation | Min | Max | Observations |
|--------------|-------------|--------|--------------------|------|-------|--------------|
| PUB TIME(MM) | 72.33866667 | 69 | 44.07649663 | 0.1 | 144 | 150 |
| FOLLORS(K) | 20212.15333 | 14300 | 18023.87129 | 28 | 63800 | 150 |
| LENGTH | 239.6066667 | 234 | 43.70660213 | 169 | 500 | 150 |
| SPEED | 4.886666667 | 5 | 1.173171384 | 2 | 7 | 150 |
| Contrast | 5.006666667 | 5 | 1.392678201 | 1 | 9 | 150 |
| beat | 5.233333333 | 5 | 1.543188831 | 2 | 8 | 150 |
| Intensity | 133.1333333 | 125 | 70.46020026 | 12 | 343 | 150 |
| elegance | 126.2 | 106.5 | 73.35845614 | 27 | 384 | 150 |
| Views(M) | 1145.754867 | 873 | 965.1479594 | 0.03 | 5345 | 150 |

3.4. Coefficient estimation

Since the model provided is linear and each variable is multiplied by a coefficient, we can estimate the characteristic of each coefficient.

Assumption 1: Coefficient for Pub Time shall be negative

The longer the song is published, the fewer people are curious about the song since they have listened to the song several times. As a result, they will listen to the song less often. The popularity will soon decrease, and the coefficient shall be negative.

Assumption 2: Coefficient for Followers shall be positive and small

Since YouTube will push every song to you published by the YouTuber you followed, the more followers a YouTuber has, the more views the video will have. Since the followers index is calculated in K, the number is much bigger than the popularity number. As a result, the coefficient must be pretty small.

Length, intensity, elegance, and extremity are determined by the listener's preference, which is impossible to estimate without linear regression.

4. Analysis

In the regression, we need to choose a significant level to do the hypothesis test. Therefore I choose 25%, 10%, 5% and 1%, symbolized *, **, ***, ****. Since the estimated data are somehow not that precise, there is a 25% level which can also reflect a new significance level.

For example, there is a hypothesis H0 that the co-relation of Popularity and Publish time is no co-relation. H1 is the alternate hypothesis which is the co-relation of Popularity and Publish time is not 0, the two are co-related.

Since the P-value of Pub Time is 3.7421E-05 smaller than 0.01, we mark it ***.

Table 2 Example for Hypothesis Test

| | Popularity vs Pub Time |
|----|------------------------|
| H0 | $\beta=0$ |
| H1 | $\beta \neq 0$ |
| r | -0.272 **** (0.06) |

In this case, the P-value is quite low and it is a **** Hypothesis and H0 is false. As a result, Popularity and Pub Time are co-related.

Since the first three variables (Pub Time, Followers, Length) are collected with precise data and show the external social effect, Model 1 is a regression based on these three variables. The other three (Intensity, Elegance, Extremity) describe the internal characteristics of the song. As a result, Model 2 is a regression based on these three variables. Model 3 combines Model 1 and Model 2, showing both external and internal effect to popularity.

4.1. Regression model

$$Y_i = 70.02 - 0.272X_1 + 4.24 \times 10^{-4}X_2 - 0.0258X_3 - 0.0914X_4 - 0.175X_5 + 6.392X_6 + \epsilon$$

Table3: Multiple Linear Regression Results and significant level

| Dependent Variable | Popularity | Popularity | Popularity |
|----------------------|---------------------------|--------------------|---------------------------|
| Independent Variable | Model1 | Model2 | Model3 |
| Publish Time | -0.250 **** (.063) | | -0.272 **** (.064) |
| Followers | 4.93E-4 **** (1.50E-4) | | 4.24E-4 **** (1.52E-4) |
| Length | -0.038 (.063) | | -0.0258 (.065) |
| Intensity | | -0.081 (.125) | -0.0914 (.115) |
| Elegance | | -0.169 * (.122) | -0.175 * (.111) |
| Extremity | | 6.14 (5.90) | 6.392 * (5.46) |
| Intercept | 40.80 *** (15.43) | 50.71 * (28.23) | 70.02 *** (31.27) |
| R Square | 0.155 | 0.039 | 0.226 |
| Adjusted R Square | 0.137 | 0.019 | 0.193 |
| Observations | 150 | 150 | 150 |

*, **, ***, **** stand for significance at the 25%, 10%, 5%, 1% levels. Standard errors are in the parentheses.

4.2. Assumption confirmation

As mentioned in the Coefficient Estimation part, the coefficient for Publish Time shall be negative. From Model1 and Model3, where Publish Time appears, we can confirm Assumption1.

Also, from Model 1 and Model3, where Followers appears, we can see that the coefficients are both positive and are to the -fourth power of 10, which is relatively small relatively. As a result, Assumption2 proves true.

4.3. Coefficient

According to the model:

As the published time increases by a month, while other conditions hold the same, views per month will decrease by 0.272 million.

As the followers of the YouTuber increase by a thousand, while other conditions hold the same, views per month will increase by 424.

As the length of the song increases one second, while other conditions hold the same, views per month will decrease by 0.0258 million.

As the intensity index increases by one, while other conditions hold the same, views per month will decrease by 0.0914 million.

As the elegance index increases by one, while other conditions hold the same, views per month will decrease by 0.175 million.

As the Extremity increases by one, while other conditions hold the same, views per month will increase by 6.392 million.

4.4. Analysis

Noted that in Model 2 and Model 3, Length and Intensity Index have pretty high p-values. As a result, these two indexes might not have a strong relationship with popularity.

The calculated the Pearson correlation coefficient (PCC) of Length vs Popularity is -0.14 and the Pearson Correlation coefficient of Intensity vs Popularity is 0.16, which are not quite high and confirm the high p-value in the model.

In Model 1 and Model 3, Published Time and Followers have ****, meaning that they have a strong relationship with popularity.

Calculate the Pearson correlation coefficient (PCC) of Pub Time vs Popularity is -0.36 and Pearson Correlation coefficient of Followers vs Popularity is 0.30, which are relatively high and confirm the low p-value in the model.

4.5. R square

R square measures how precisely the model can fit every point. When R square equals 1, the model can fit every point on the graph accurately. R square is between 0 and 1.

However, as the variables increase and R square increases. We need to consider the number of variables.

Then we need to calculate the adjusted R square, which takes the number of variables into consideration.

In Model 2, both R square and adjusted R square are pretty low. However, all of the three variables in Model 2 are calculated by Beat, Contrast, and Speed Index, estimated by me. Basically, errors are acceptable.

What's more, we can see that the Adjusted R Square in Model 3 is larger than that in either Model 1 or Model 2. Since Model 3 has three more variables than both Model 1 and Model 2, also Adjusted R square increases, we can infer that both internal characteristics (Intensity, Elegance, and Extremity) and external data (published time, followers, and length) do have positive influences to fix the model.

Proof

$$R^2 = \frac{SST - SSE}{SST}$$

let

$$\hat{y}^k = a + b_1 * x_1 + b_2 * x_2 + \dots + b_k * x_k$$

$$y = \hat{y}^k + u$$

$$SSE_k = \sum_{i=0}^n (y_i - \hat{y}_i^k)^2$$

We also

$$\hat{y}^{k+z} = a + b_1 x_1 + b_2 x_2 + \dots + b_k x_k + \beta z$$

$$y = \hat{y}^{k+z} + v$$

$$\begin{aligned} SSE_{k+z} &= \sum_{i=0}^n (y_i - \hat{y}_i^{k+z})^2 = \sum_{i=0}^n \{(y_i - \hat{y}_i^k) - \beta z_i\}^2 \\ &= SSE_k + \sum_{i=0}^n \{(\beta \times z_i)^2 - 2(y_i - \hat{y}_i^k)(\beta \times z_i)\} \end{aligned}$$

Since the can change the value of $b_1, b_2, b_3, \dots, \beta$ and have the SSE minimized, Also, β can always be zero, while $SSE_{k+z} = SSE_k$

Otherwise, $SSE_{k+z} < SSE_k$ As a result, $R^2_{k+z} \geq R^2_k$

Conclusion: As the number of variables increases, R square increases.

4.6. Optimal case

Interestingly, the coefficients for Intensity and Elegance are all negative, while the coefficient for Extremity is positive. It is quite confusing since high Extremity must be a result of either high intensity or high elegance. The proper explanation might be that the coefficients for Intensity and Elegance are negative to fix the high positive coefficient for Extremity or other characteristics. In order to check the optimal case, we can use a computer algorithm.

Since the Published Time, Followers, and length are in a linear relation to popularity, in the optimal case, we shall just make the published time lower, followers higher, and length shorter. We cannot take these three variables into consideration in the code.

We consider beat(b), contrast(c), and speed(s) in the code since the Intensity, Elegance, and Extremity are all variations of these three original data. Setting these three variables to positive integers from 1 to 10, we can get the result that the optimal case is Speed 9, Beat 9 and Contrast 2 with the distribution of -6.56 to the popularity. The lowest contribution takes place when the three are all equal to 1. Since the three variables are reciprocal, the exchange of the three index values does not change the total contribution to popularity. We have three different situations where we can get the same contribution to popularity: Speed 9, Beat 9, Contrast 2; Speed 9, Beat 2, Contrast 9; Speed 2, Beat 9, Contrast 9. In the first case we can infer that pop songs with a fast speed and strong sense of beat thoroughly are quite popular. In the second and third cases we can infer that if the sense of beat is weak or the speed is pretty slow, with a strong contrast, the songs can still be popular.

However, since we can hardly have an index graded 1 or 10, we adjust the range to Speed [2,7], Beat[2,8], and Contrast[1,9] (highest and lowest grade in every index). The optimal case is Speed 2 Beat 8 and Contrast 9, with a contribution of -9.31.

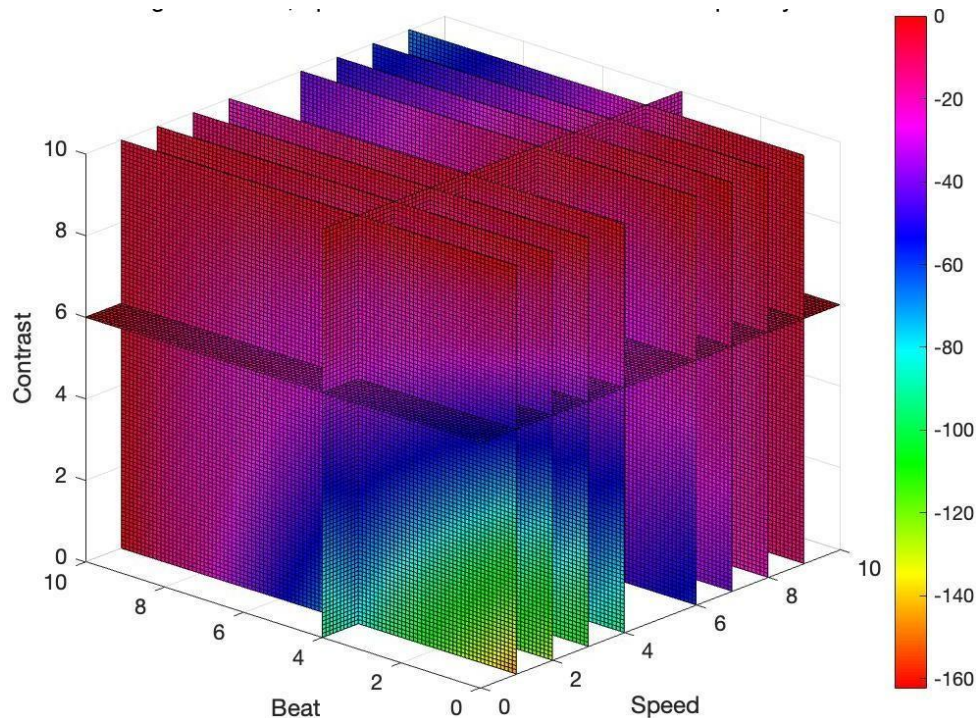


Figure 13: Beat, Speed and Contrast Contribution to Popularity

In Figure 13, we can see that the purple part is the highest contribution to popularity while the orange part is the lowest. Most songs with high contrast have much higher popularity. Those which are slow, elegant, immutable are not so welcomed in the current market.

In Appendix B, the code is for the case where all s , b , c are not integers. We determined Speed 7 Beat 8 and Contrast 3.69. We can also change the function in the algorithm to give different weights to each index and check other cases using this algorithm.

5. Conclusion

In this article, the regression model of a song's popularity is proposed. From the coefficient, we can observe that the increases of followers and extremity do have positive effects on a song's popularity, while those of Published time, length, intensity, and elegance do have negative effects.

In such a competitive pop song market, consumers prefer songs with strong extremity, no matter whether it is intense or elegant. Composers can try to make the music a strong contrast between different sections—for example, a slow beginning with an exciting chorus. While composing and making music, try to avoid making the music peaceful thoroughly.

However, Followers Index is still a significant index among all. One thousand more followers can lead to 424 more views per month. As a result, those YouTubers or singers with more than 50 million followers will have about 21 million more views per month than those who recently entered the music market.

More work could be done by adjusting and trying different models, such as the logarithm model. These can be done by varying the regression option and re-write the function in the Matlab code.

References

- [1] Sloan, *CONTEMPORARY MUSIC REVIEW*, Taylor Swift and the Work of Songwriting, 2021. <https://www.tandfonline.com/doi/full/10.1080/07494467.2021.1945226>
- [2] Bernardo Monechi, Pietro Gravino, Vito D. P. Servedio, Francesca Tria and Vittorio Loreto, *The Royal Society, Significance and popularity in music production*, 2017. <https://royalsocietypublishing.org/doi/10.1098/rsos.170433#RSOS170433F3>
- [3] David Martín-Gutiérrez; Gustavo Hernández Peñaloza; Alberto Belmonte-Hernández, *IEEE, A*

Multimodal End-to-End Deep Learning Architecture for Music Popularity Prediction, 2020. <https://ieeexplore.ieee.org/document/9007339>

[4] Sciandra, Mariangela; *Universita degli Studi di Palermo Dipartimento di Scienze Economiche Aziendali e Statistiche*, Spera, Irene Carola; *University of Palermo, Journal of Applied Statistics*, A model-based approach to Spotify data analysis: a Beta GLMM, 2020. <https://iris.unipa.it/retrieve/handle/10447/429276/939121/JAS.pdf>

[5] Martin-Gutierrez, D (Martin-Gutierrez, David) ,Penaloza, GH (Hernandez Penaloza, Gustavo) ,Belmonte-Hernandez, A (Belmonte-Hernandez, Alberto) ,Garcia, FA (Alvarez Garcia, Federico), *IEEE Access*, A Multimodal End-to-End Deep Learning Architecture for Music Popularity Prediction, 2020. <https://ieeexplore.ieee.org/document/9007339>