

APPLICATION NOTE



# A model-based approach to Spotify data analysis: a Beta GLMM

Mariangela Sciandra<sup>a</sup> and Irene Carola Spera<sup>b</sup>

<sup>a</sup>Department of Economics, Business and Statistics, University of Palermo, Palermo, Italy; <sup>b</sup>University of Palermo, Palermo, Italy

## ABSTRACT

Digital music distribution is increasingly powered by automated mechanisms that continuously capture, sort and analyze large amounts of Web-based data. This paper deals with the management of songs audio features from a statistical point of view. In particular, it explores the data catching mechanisms enabled by Spotify Web API and suggests statistical tools for the analysis of these data. Special attention is devoted to songs popularity and a Beta model, including random effects, is proposed in order to give the first answer to questions like: which are the determinants of popularity? The identification of a model able to describe this relationship, the determination within the set of characteristics of those considered most important in making a song popular is a very interesting topic for those who aim to predict the success of new products.

## ARTICLE HISTORY

Received 19 November 2019  
Accepted 26 July 2020

## KEYWORDS

Spotify web API; audio features; popularity index; Beta GLMM

## 2010 MATHEMATICS SUBJECT CLASSIFICATIONS

62; 62H; 62P

## 1. Introduction

Music plays an important role in everyday life of people, and with digitalization, large collections of musical data are formed, which tend to be further cumulated by music lovers [23]. This has led to music collections, not only on the private shelf as audio or video discs and domain discs, but also on the hard disk and online, to grow beyond what was previously impossible. With the advent of new technologies, it has become impossible for a single individual to keep track of the music and the relationships between different songs. The techniques of data mining and automatic learning can help the navigation in the world of music [14].

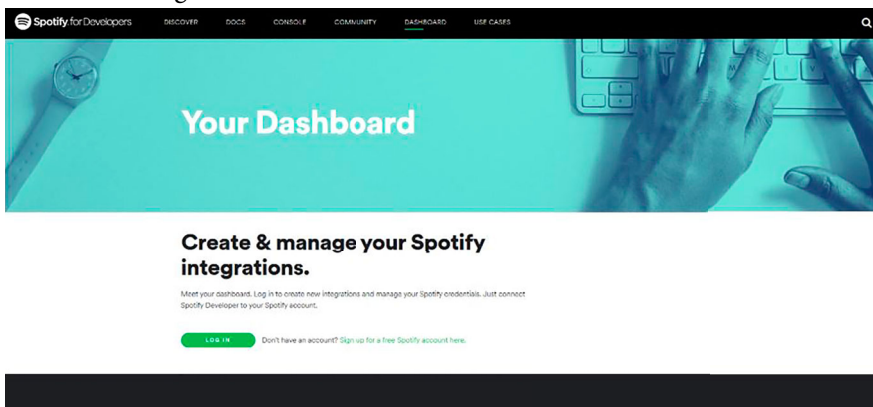
Data mining strategies are often based on two main problems: the type of available data and the use you want to make of them. What kind of data is the music? A collection of music tracks consists of various types of data; for example, data could consist of music audio files or metadata such as track title and artist name [20]. What kind of analysis can be carried out? The musical data mining provides specific methods to answer to the most varied questions: e.g. gender classification, identification of artists/singers, mood/emotion detection, instrument recognition, similarity search music, musical synthesis and so on.

This research investigates the relationship between song data audio features obtained from the Spotify database (e.g. key and tempo) and song popularity, measured by the number of streams that a song has on Spotify. Previous researches on the topic of new product success prediction have identified multiple approaches to answer to this question. Moreover, the existing body of research, which defines many popularity prediction models stresses the complexity of the mechanisms of song popularity. Research from [13] shows that it is feasible to predict the popularity metrics of a song significantly better than random chance based on its audio signal. Additionally, Ni *et al.* [18] also show that certain audio features such as Loudness, duration and harmonic simplicity correlate with the evolution of musical trends. Dhanaraj and Logan [7] propose features from both songs' lyrics and audio content for prediction of hits and also study a hit detection model based solely on lyrics' features. In an attempt to predict the popularity of a song from Spotify's song data, the research of Berger [2] uses (Echo-Nest) audio features similar to this research and uses Spotify's own calculated metric 'popularity' to measure popularity. Other attempts through classical linear regression or quadratic models can be found on the net but they are not exhaustive works and do not take into account the particular data structure and other aspects that could lead to biased predictions. This is the reason why this paper can be considered as an innovative way to look at popularity predictions and represents an innovative approach inside the literature. Paper is organized in the following way: in Section 2, the way to connect Spotify Web API to the R software is explained; available song's Spotify audio features are described in Section 3; Section 4 is devoted to the proposal of a new class of models to predict song's popularity as a function of the Spotify audio features. An application to a real dataset is carried out in Section 5, future work and conclusions follow.

## 2. Spotify and R

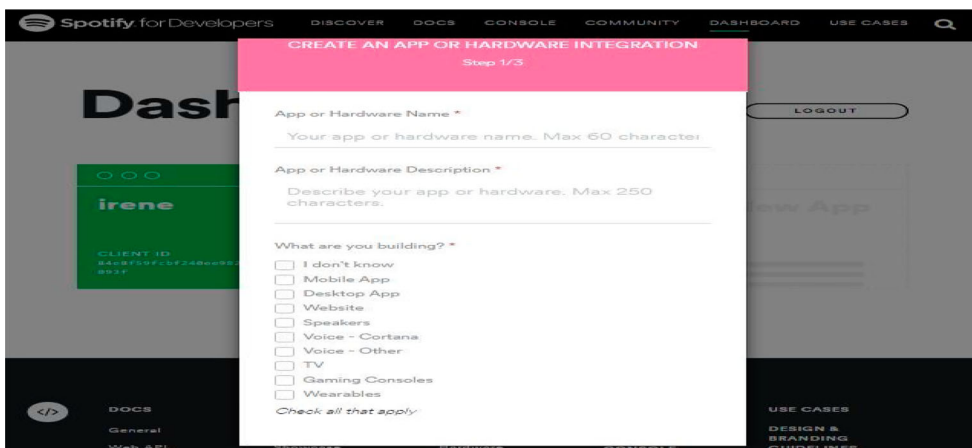
Spotify [24] is one of the most famous music Apps in the world, a program of the last generation grown up exponentially over the last few years. This platform allows you to listen to music streaming to computers, smartphones and tablets, choosing from more than 30 million tracks, old and new, of the main international record companies, without having to purchase individual songs or albums legally.

In this section, a brief description of the steps necessary to connect Spotify Web API to the R software are reported. First, you have to create an app on the Spotify's developer platform, accessing this dashboard:

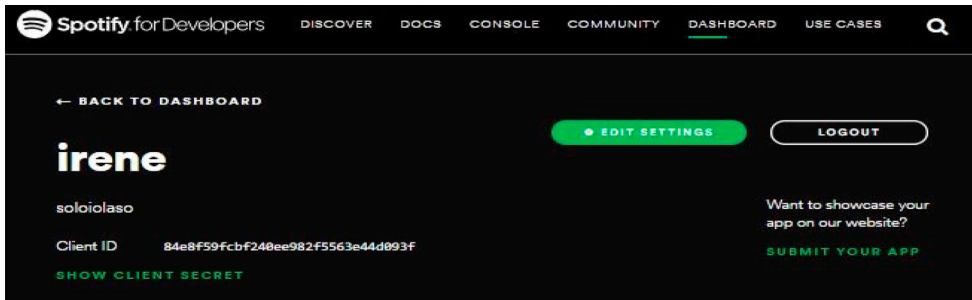


After accessing the page, you have to login and create an app:

- (1) give a name to the application and provide a brief description;
- (2) then you have to specify the purpose for which it is created (example for commercial purposes, teaching, etc. . . .);
- (3) finally, accept the conditions to complete the account configuration.



After this procedure, user identification code and password are generated, both are alphanumeric codes that are essential to release the access token to connect Spotify with R.



The packages used to make the connection are:

- *rvest*: it allows you to extract data from a web page (web scraping);
- *tidyverse*: it is used for data transformation and cleaning;
- *DSpoty*: it extracts the song's audio features from Spotify's Web API.

After loading the R libraries listed above, the following functions are used to connect and extract data from Spotify [5]:

```
Sys.setenv(client_id = 'xxxxxxxxxxxxxxxxxxxxxxxx')
Sys.setenv(client_secret = 'xxxxxxxxxxxxxxxxxxxxxxxx')
access_token ← DSpoty::get_spotify_access_token()
Liga ← get_artist_audio_features('Ligabue')
```

They are especially useful for:

- `Sys.setenv()`: storing in R the identification code and password to access Spotify;
- `get_spotify_access_token()`: generating the token to access all data using the Spotify library;
- `get_artist_audio_features()`: extracting all the audio characteristics of the songs that make up the discography, by specifying the name of the singer.

### 3. Spotify audio features

Spotify Web API makes users able to extract several audio features of songs. The available features are listed in Table 1.

The aim of this paper is to investigate if audio features from Spotify can be considered as determinants of the stream popularity of songs [19].

#### 3.1. Track API popularity

Among all the features returned by Spotify, song popularity plays an important role. The popularity of a track is a value between 0 and 100, with 100 the most popular. The popularity is calculated by algorithm and is based, in the most part, on the total number of track plays and taking into account how recent those plays are. Generally speaking, songs that are being played a lot now will have higher popularity than songs that were played a lot in the past. Duplicate tracks (e.g. the same track from a single and an album) are rated independently. Artist and album popularity is derived mathematically from track popularity. Note that the popularity value may lag actual popularity by a few days: the value is not updated in real-time.

Songs popularity is an important issue for the music industry. In 2017, the music industry generated \$8.72 billion in the United States alone. Thanks to growing streaming services (Spotify, Apple Music, etc.) the industry continues to flourish. The top 10 artists in 2016 generated a combined \$362.5 million in revenue. The question of what makes a song popular has been studied before with varying degrees of success [9]. Every song has key characteristics including lyrics, duration, artist information, temp, beat, Loudness, chord, etc. Previous studies that considered lyrics to predict a song's popularity had limited success.

### 4. Using Spotify data to predict what songs will be hits

The aim of this section is the identification of the determinants of songs' popularity. In particular, we want to investigate the possible relationship between the audio characteristics of the songs in the Spotify database (for example, Energy, Loudness, etc. . . .) and the popularity of the songs also available in the Spotify dataset. The identification of a model able to describe this relationship, the determination within the set of characteristics of those considered most important in making a song popular is a very interesting topic for those who aim to predict the success of new products. Then, the fundamental question is: What does determine popularity? Why is a song popular?

In cultural markets like music, forecasting is very complex. Studies in this field called Hit Song Science (HSS) are of interest to record companies but also to consumers themselves and to Spotify [16]. Previous attempts in this direction have always referred to linear

**Table 1.** Spotify audio features.

acousticness	float	A confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic.
analysis_url	string	An HTTP URL to access the full audio analysis of this track. An access token is required to access this data.
danceability	float	Danceability describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. A value of 0.0 is least danceable and 1.0 is most danceable.
duration_ms	int	The duration of the track in milliseconds.
energy	float	Energy is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity. Typically, energetic tracks feel fast, loud, and noisy. For example, death metal has high energy, while a Bach prelude scores low on the scale. Perceptual features contributing to this attribute include dynamic range, perceived loudness, timbre, onset rate, and general entropy.
id	string	The Spotify ID for the track.
instrumentalness	float	Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Rap or spoken word tracks are clearly "vocal". The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content. Values above 0.5 are intended to represent instrumental tracks, but confidence is higher as the value approaches 1.0.
key	int	The key the track is in. Integers map to pitches using standard <a href="#">Pitch Class notation</a> . E.g. 0 = C, 1 = C#/Db, 2 = D, and so on.
liveness	float	Detects the presence of an audience in the recording. Higher liveness values represent an increased probability that the track was performed live. A value above 0.8 provides strong likelihood that the track is live.
loudness	float	The overall loudness of a track in decibels (dB). Loudness values are averaged across the entire track and are useful for comparing relative loudness of tracks. Loudness is the quality of a sound that is the primary psychological correlate of physical strength (amplitude). Values typical range between -60 and 0 db.
mode	int	Mode indicates the modality (major or minor) of a track, the type of scale from which its melodic content is derived. Major is represented by 1 and minor is 0.
speechiness	float	Speechiness detects the presence of spoken words in a track. The more exclusively speech-like the recording (e.g. talk show, audio book, poetry), the closer to 1.0 the attribute value. Values above 0.66 describe tracks that are probably made entirely of spoken words. Values between 0.33 and 0.66 describe tracks that may contain both music and speech, either in sections or layered, including such cases as rap music. Values below 0.33 most likely represent music and other non-speech-like tracks.
tempo	float	The overall estimated tempo of a track in beats per minute (BPM). In musical terminology, tempo is the speed or pace of a given piece and derives directly from the average beat duration.
time_signature	int	An estimated overall time signature of a track. The time signature (meter) is a notational convention to specify how many beats are in each bar (or measure).
track_href	string	A link to the Web API endpoint providing full details of the track.
type	string	The object type: "audio_features"
uri	string	The Spotify URI for the track.
valence	float	A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track. Tracks with high valence sound more positive (e.g. happy, cheerful, euphoric), while tracks with low valence sound more negative (e.g. sad, depressed, angry).

or quadratic model regression [17]. In this paper, the application of a Beta regression with random effects is proposed. The choice of this class of models derives by nature of the response variable (a continuous variable limited in  $[0, 100]$ ) and by the correlation structure in the data: it is assumed, in fact, that songs belonging to the same album can be related to each other more than song from different albums; ignoring this level of hierarchy in the data could lead to biased or inefficient results.

**4.1. Beta regression for correlated data**

In this work, we study the dependence of the popularity of songs on the musical characteristics by using an extension of the Beta regression model, including random effects. The resulting model will be a generalized Beta model with mixed effects (Beta GLMM). Before defining the model from a theoretical point of view, following a brief remained to the classical Beta regression and the Generalized Linear models with Mixed Effects (GLMMs); this brief summary will be useful to understand the reasons why we have chosen to focus on this particular model and for the theoretical definition of the model itself.

**4.1.1. Beta regression**

Beta distribution is a continuous probability distribution defined in the unitary range with a density function given by

$$f(y, \mu, \phi) = \frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1 - \mu)\phi)} y^{\mu\phi-1} (1 - y)^{(1-\mu)\phi-1}, \tag{1}$$

where  $\Gamma(\cdot)$  indicates the Gamma function. The parameter  $\mu$  indicates the expected value of  $Y$ , i.e.  $E(Y) = \mu$ . The parameter  $\phi$  meets the definition of a precision parameter because, for fixed  $\mu$ , the higher the value of  $\phi$ , the lower the variance of the dependent variable. More specifically,

$$\text{Var}(Y) = \frac{\mu(1 - \mu)}{1 + \phi}. \tag{2}$$

In Beta regression models [8], the parameter that indicates the average  $\mu \in (0, 1)$  of the Beta distribution is expressed as a function of the covariates, while the parameter of precision  $\phi \in R^+$  is treated as a disturbance parameter. In order to ensure that the linear predictor takes on values in the space given by the dependent variable’s support, the link logit represents the most commonly chosen link function

$$g(\mu_i) = \log\left(\frac{\mu_i}{1 - \mu_i}\right) = x_i^T \beta, \tag{3}$$

where  $x_{ij}^T$  denotes a vector of explanatory variables, and  $\beta$  refers to the vector of regression coefficients,  $i = 1, \dots, N$ . The Beta distribution is defined only on the open unit interval. If exact one and zero values are admitted, these values must be transformed in order to ensure the nature of the Beta distribution support [3]. The most frequently applied transformation is:

$$Y^* = [Y(N - 1) + 0.5]/N \tag{4}$$

where  $Y^*$  is the transformed and  $Y$  is the untransformed dependent variable. Alternatively, it was suggested to add a small amount of  $\varepsilon$ , for example 0.005 or 0.01 to the lower limit, and



subtract the same amount from the upper limit. Hunger et al. [11] also observe that when the resulting values are too close to the boundary points, the accuracy of the estimates may decrease significantly.

**4.1.2. Generalized linear mixed models (GLMM)**

Generalized Linear Mixed Models (or GLMM) are an extension of the Generalized Linear Model (GLM) in which the linear predictor contains random effects in addition to the usual fixed effects. For this model class, the assumption of homogeneity and independence of the sample units is lost. In addition, with regard to the distribution of the response variable, the GLMM inherit from the GLM the idea of extending mixed linear models to the non-normal data case [15]. GLMM provide a wide range of models for the analysis of data that have some form of grouping, since differences between groups can be modeled through the use of a random effect. The basic concept is the structure in *cluster*: the data with clustering structure has a univariate response variable  $y$  double indexed,  $i$  for the first level units and  $j$  for the second level units and a vector  $x_{ij}$  of explanatory variables  $p$  for the  $j$ th unit in the  $i$ th cluster. It is important to remember that clusters can have different sizes and that this can influence the results of the analysis.

These models are useful in the analysis of many types of data, including longitudinal data. The general form of the model, in matrix notation, is

$$y = X\beta + Zb + \varepsilon, \tag{5}$$

where  $y$  is a column vector  $N \times 1$ ;  $X$  is an array  $N \times p$  of explanatory variable  $p$ ;  $\beta$  is a column vector  $p \times 1$  of fixed effect regression coefficients;  $Z$  is the random effects model matrix  $N \times q$  for random effects  $q$ ;  $b$  is a vector  $q \times 1$  of random effects;  $\varepsilon$  is the column vector  $N \times 1$  of residues. The assumptions that underlie this class of models can be summarized as follows:

$$\begin{aligned} Y_{ij} \mid (x_{ij}, z_{ij}, b_i) &\sim C.\xi.N(\theta_{ij}, \phi); \\ \eta_{ij} &= X_{ij}^T \beta + Z_{ij}^T b_i; \\ g(\mu_{ij}) &= \eta_{ij}; \\ \mu_{ij} &= h(\eta_{ij}) = g(\eta_{ij})^{-1}; \\ b_i &\sim (0, \Sigma_q); \\ Y_i &\perp Y_j \quad \forall i \neq j. \end{aligned}$$

By putting together all the assumptions the conditional distribution is easily derived

$$f(y_{ij} \mid x_{ij}, z_{ij}, b_i) = \exp \left\{ \frac{y_{ij}\theta_{ij} - b(\theta_{ij})}{\phi} + c(y_{ij} + \phi) \right\}.$$

Conditioning on  $b_i$ , observations from the same cluster are assumed to be independent. In addition, the conditional expected value is related to the linear predictor (containing both random and fixed effects) by the following linking function  $g(\cdot)$ :

$$g(\mu_{ij}) = x_{ij}^T \beta + Z_{ij}^T b_i.$$

### 4.1.3. The Beta GLMM

In longitudinal analyses or when subjects have any grouping structure, observations related to the same unit will typically be correlated, violating the assumption of independence of observations typical in regression models. The dependence within clusters can be accounted for by adding random cluster or subject effects in the linear predictor [3]. Consider the case of longitudinal studies where  $j = 1, \dots, n_i$  observations are nested within  $i = 1, \dots, N$  subjects. Let  $b_i$  denote the vector of random effects specific to each subject  $i$ .

Adding random effects to the beta regression model in (3.3), we get the GLMM beta [3] given by

$$\log\left(\frac{\mu_{ij}}{1 - \mu_{ij}}\right) = x_{ij}^T \beta + z_{ij}^T b_i \quad \text{con } b_i \sim N(0, G), \tag{6}$$

where  $z_{ij}^T$  is a vector of explanatory variables and  $G$  is the defined positive covariance matrix of random effects. Note that although the assumption of normality for random effects is common and statistically convenient, other distribution hypotheses are also possible. In a longitudinal study, the  $b_i$  is typically a scalar (for random intercept models) or a bivariate vector (for random intercept models with a random regression coefficient), i.e.  $z_{ij}^T = (1, t_{ij})$ , where  $t_{ij}$  is the measurement time  $j$  for the subject  $i$ .

In the Beta GLMM, the regression parameters have only one specific interpretation per unit and do not describe the effect of the respective variable on the population average; this is due to the non-linear transformation of the average response (i.e. the logit link) as it can be deduced that

$$\text{logit}(E(Y_{ij} | b_i)) = x_{ij}^T \beta + z_{ij}^T b_i, \tag{7}$$

but

$$\text{logit}(E(Y_{ij} | b_i)) \neq x_{ij}^T \beta.$$

Model parameters can be estimated by maximizing the marginal probability that is obtained by integrating the joint distribution of  $[Y, \mathbf{b}]$  on random effects. The contribution to the log-likelihood by each group is as follows:

$$f_i(y_i | \beta, \Sigma, \phi) = \int \prod_{j=1}^{n_i} f_i(y_{ij} | b_i, \beta, \phi)(b_i | \Sigma) db_i. \tag{8}$$

Assuming independence among the  $N$  groups, the full likelihood is

$$L(\beta, \Sigma, \phi) = \prod_{i=1}^N f_i(y_i | \beta, \Sigma, \phi). \tag{9}$$

## 5. An application: Luciano Ligabue’s audio Spotify features analysis

In order to apply the proposed model to a real case, the whole discography of Luciano Ligabue has been analyzed. Luciano Ligabue is one of the most successful Italian artists. He is a famous singer-songwriter, film director and writer. In over 30 years of his career, he has won more than 60 awards for his musical activity, 5 awards for his work as a writer and





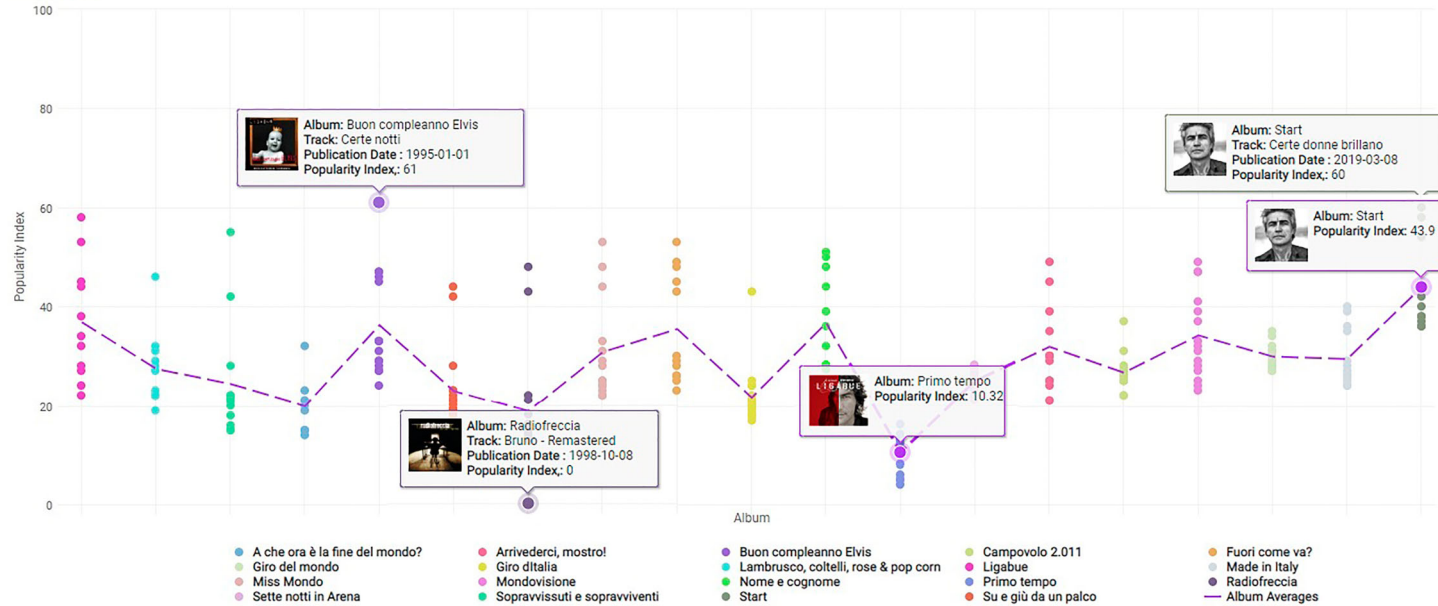
**Figure 1.** Average values of Spotify features over time.

12 awards for his film work [21]. In this application 19 albums for a total of 273 songs have been considered. Among the information available on Spotify there is also the year in which the song was published. It is clear that, statistically speaking, the year represents a proxy for the album because there are no different albums released in the same year. So, as the album has been assumed as a grouping variable, the year has not been considered among the covariates. However, for completeness of information, the graph below (Figure 1) shows the average values over time of the examined characteristics.

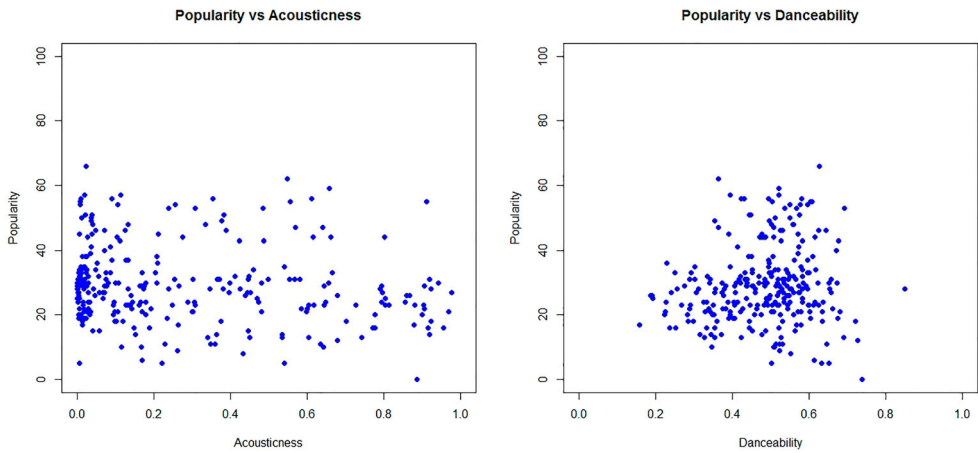
In terms of popularity, it follows that, after reaching the bottom point of popularity compared to the entire career, in 2007, there is an increase in popularity which reaches higher values ever until it reaches its absolute maximum in 2019 with ‘Start’. The musical characteristics do not show any particular trend; it is as if the musical genre is constant over time, except for the peaks of Speechiness of 1998 given by ‘Radiofreccia’. It is important to note that the peak in Speechiness is accompanied by a corresponding reduction in Energy, Danceability and Loudness. In addition, Speechiness seems to have a turning point in 1998: the Speechiness index for the years prior to 1998 is always higher than the Speechiness index of the subsequent years. From the plot, it is possible to distinguish between all the live albums, in which the values of the homonymous index reach abnormal peaks, and then assume values around 0.2 in the case of non-live album. Among all the features, the Loudness and the Instrumentalness of the songs seem to be those more stable within the entire discography of Luciano Ligabue.

Focusing on the popularity over the years, the plot in Figure 2 shows that, on average, the most popular album is *Start*, released on 8 March 2019, followed by the album *Buon Compleanno Elvis* in 1995, while the least listened to album is *Primo tempo*, probably because it is a collection. In particular, the most-streamed songs are ‘Polvere di stelle’ and ‘Certe Notti’, while the less popular song is ‘Radiofreccia’, soundtrack of the movie with the same name. This song is only Instrumentalness and not spoken, so it could have an influence on its low popularity.

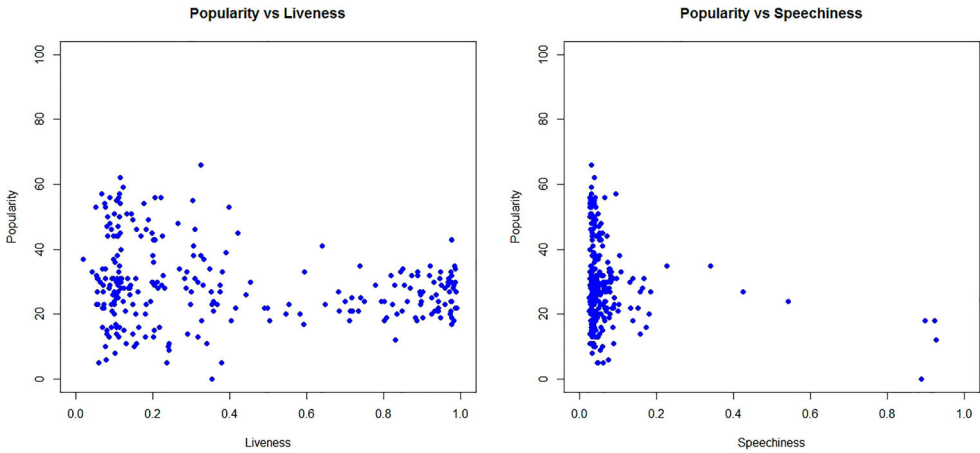
## Discography



**Figure 2.** Distribution of songs according to the Popularity Index, conditioning to the album (13 November 2019).



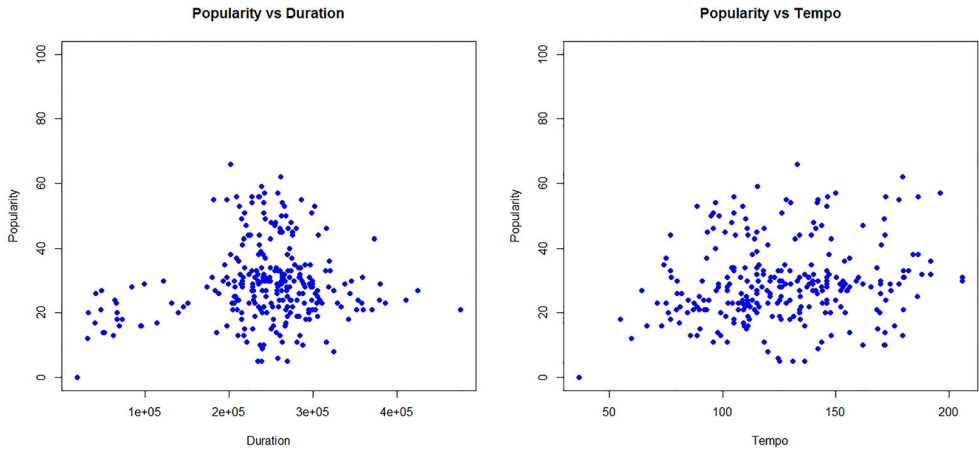
**Figure 3.** Scatterplot of popularity vs. acousticness and danceability.



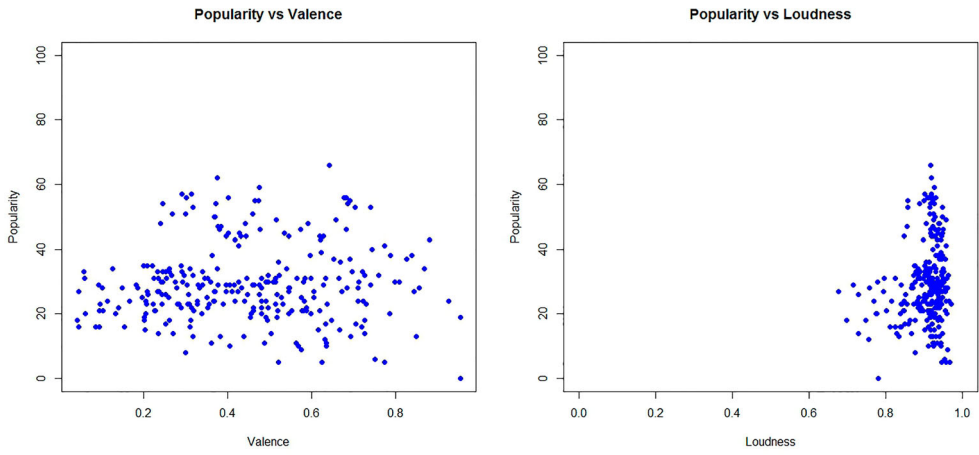
**Figure 4.** Scatterplot popularity vs. liveness and speechness.

In the light of trends observed in Figure 1, before fitting the model we want to investigate if there are correlations among these variables in order to avoid problems of multicollinearity. Scatterplots in which the individual audio characteristics are related with the popularity index of the songs are shown in Figures 3 and 4.

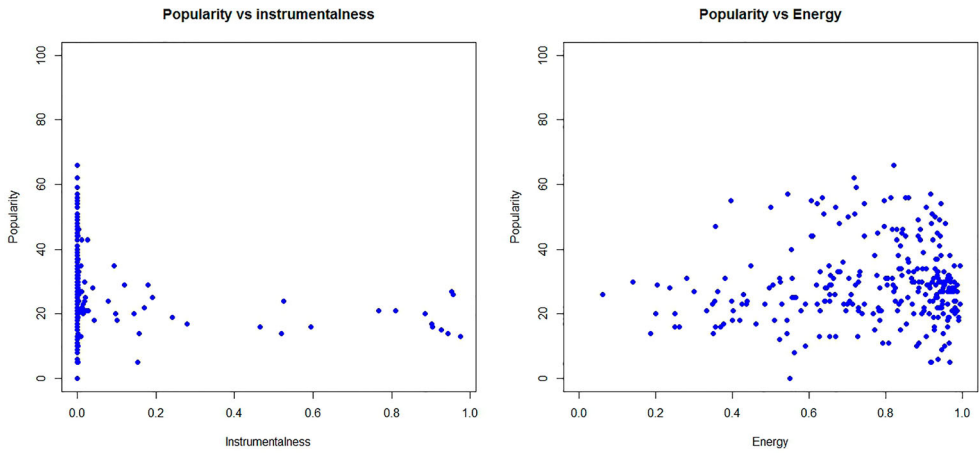
Scatterplots show that most of the songs in Luciano Ligabue's discography are not very Speechness, not very Instrumentalness, with a lot of Energy and have Loudness values close to zero. In particular, there is a downward trend in Popularity compared to Live, Speechness, Instrumentalness and Acoustics, i.e. as these audio characteristics increase, Popularity decreases. The Valence, Rhythm and Danceability charts do not show any particular trend, while Energy seems to have a positive trend, i.e. Energy increases and the Popularity of the song increases as well (Figures 5–7).



**Figure 5.** Scatterplot popularity vs. time and duration.



**Figure 6.** Scatterplot popularity vs. valence and loudness.



**Figure 7.** Scatterplot popularity vs. instrumentality and energy.

**Table 2.** Selected Beta GLMM.

Random effects					
Groups	Variance	Std. Dev.			
Album name (Intercept)	0.1926	0.4389			
Fixed effects					
	Coeff.	Std. Error	z value	Pr(> Z )	
<i>Intercept</i>	−1.18120	0.21596	−5.470	4.51e−08	***
<i>Speechness</i>	−1.52742	0.33042	−4.623	3.79e−06	***
<i>Instrumentalness</i>	−0.34554	0.18122	−1.907	0.05655	.
<i>Liveness</i>	−0.27380	0.12229	−2.239	0.02516	*
<i>Duration</i>	0.08419	0.03061	2.751	0.00594	**
<i>Energy</i>	0.66948	0.26992	2.480	0.01313	*
<i>Valence</i>	1.16800	0.47151	2.477	0.01324	*
<i>Energy:Valence</i>	−1.46232	0.58209	−2.512	0.01200	*
Overdispersion parameter for the beta family				$\phi = 38.7$	

### 5.1. Model selection in Beta GLMM

In order to identify which are the musical characteristics influencing popularity, and if this relationship is similar for all albums (and therefore overtime), a *Beta model with random effects* has been estimated, since the available data has a cluster structure given by the music albums, and it has been assumed that the response variable is distributed according to a Beta distribution, as the popularity index is a continuous and limited variable in an interval the (0, 1). In addition, since to the presence of exact zeros, the index is rescaled in the bounded interval.

Through the *glmmTMB* library [4], different models were estimated and variables were selected on the basis of several information criteria. Variable selection for mixed-effects models represents a wide research topic in the literature. Some authors proposed measures for testing hypothesis on the variance components (see [6,12]) in order to detect whether an individual random component is significant or not [22]. Yet, in our application, the fitted model includes only a random intercept, so models comparison only focused on the fixed part of the model. When the random part selection is out of interest, a natural choice would be to base model selection on the *Akaike Information Criterion* (AIC) [1] from the marginal model, i.e. the model with the random effects integrated out. This leads to a biased criterion (see [10]). As an alternative we based our comparison on the conditional AIC (cAIC) introduced by Vaida and Blanchard [25], that is an extension of the classical AIC for mixed-effects models assuming the variance parameters of the random effects to be known. The final model with the lower cAIC (cAIC= −607.80) is the one shown in Table 2.

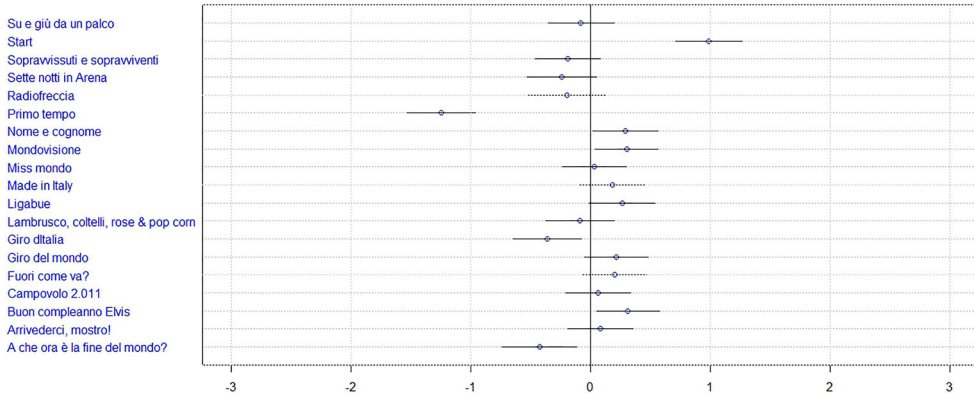
The results of the selected model show that *Speechness*, *Instrumentalness* and *Live* are the features that negatively affect the Popularity Index, while *Energy*, *Valence* and *Duration* of the song are the ones that positively affect it. It should also be noted that the interaction between *Energy* and *Valence* has a particularly negative effect on the considered index.

In order to check for overdispersion in data, a dispersion test to verify the hypothesis  $\phi > 1$  was carried out with no significant results ( $p = 0.808$ ).

Table 3 shows the changes, in percentage terms, in the Popularity Index of songs, for unitary changes of the audio features.

**Table 3.** Variations in % in popularity index.

Audio features	Variation in popularity index (%)
Speechness	-32
Instrumentalness	-9
Liveness	-7
Duration	2
Energy	16
Valence	26
Energy:valence	-31



**Figure 8.** Random intercepts with 95% confidence intervals.

**Table 4.** cAIC comparison between Beta GLMM and Normal LMM.

Random effects model	cAIC
LMM	-540.34
Beta GLMM	-607.80

Through the *Caterpillar plot* let us try to understand whether the clustered structure is relevant for obtaining better, in inferential terms, predictions.

Figure 8 confirms that the random part in the model cannot be ignored, as there is a cluster effect due to the album entitled *A che ora è la fine del mondo?*, *Buon Compleanno Elvis*, *Giro d'Italia*, *Mondovisione*, *Nome e Cognome*, *Primo tempo* and *Start*; for these albums estimate confidence intervals do not intersect the zero, this means that the hypothesis of homogeneity of the albums is rejected because the albums are different from each other. This is also confirmed by the estimate of the variance of random effects which is 0.1926 and therefore not negligible. Once the model has been selected, in order to further stress the advantage in using a Beta distribution instead of the Normal distribution so far used in literature, we compared the Beta GLMM and the Normal LMM in terms of cAIC as a measure of formal diagnostic on the distributional assumption. cAIC from both models can be found in Table 4.

## 6. Conclusions

In this work, a new class of models for dealing with songs popularity index has been introduced. The use of a Beta GLMM allows to account for clustering structure of data from music album and results on the real case example show that this structure cannot be ignored. The Spotify Web API audio features, used as covariates, have shown that not all the Spotify characteristics have high explanatory power for a higher stream count but some of them are actually important. Significant relationships were found, which lays a promising foundation for the research in prediction with these variables. In particular, Speechness, 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.

This research contributes to further understanding in the field of HSS and the new product success prediction. Creating effective prediction models is an interesting next step to this research, and so the next step would be to expand on the variables used. We hope that this paper will have practical implications also on Spotify, suggesting, for example, interesting ideas to further develop its database with the hope that data of increasing quality can lead to interesting discoveries and added value to the world of HSS.

## Disclosure statement

No potential conflict of interest was reported by the author(s).

## References

- [1] H. Akaike. *Information Theory as an Extension of the Maximum Likelihood Principle*, Second international symposium on information theory. Petrov, Boris Nikolaevich and Csaki, F, 1973, pp. 267–281.
- [2] W. Berger. Why is this song popular? (feat spotify). Available at <https://medium.com/@albert.w.berger/what-makes-a-song-popular-in-a-certain-country->
- [3] W.H. Bonat, P.J. Ribeiro, and W.M. Zeviani, *Likelihood analysis for a class of beta mixed models*, J. Appl. Stat. 42 (Aug 2014), pp. 252–266. <https://doi.org/10.1080/02664763.2014.947248>.
- [4] M.E. Brooks, K. Kristensen, K.J. van Benthem, A. Magnusson, C.W. Berg, A. Nielsen, H.J. Skaug, M. Maechler, and B.M. Bolker, *glmmTMB balances speed and flexibility among packages for zero-inflated generalized linear mixed modeling*, R. J. 9 (2017), pp. 378–400. Available at <https://journal.r-project.org/archive/2017/RJ-2017-066/index.html>.
- [5] Charlie, Rcharlie web site. Available at <https://https://www.rcharlie.com/>, 2019.
- [6] Z. Chen and D.B. Dunson, *Random effects selection in linear mixed models*, Biometrics 59 (2003), pp. 762–769.
- [7] R. Dhanaraj and B. Logan, *Automatic prediction of hit songs*, in ISMIR 2005, 6th International Conference on Music Information Retrieval, 11–15 September 2005, Proceedings, HP Laboratories Cambridge, London, UK, 2005, pp. 488–491. Available at <http://ismir2005.ismir.net/proceedings/2024.pdf>
- [8] S. Ferrari and F. Cribari-Neto, *Beta regression for modelling rates and proportions*, J. Appl. Stat. 31 (2004), pp. 799–815. <https://doi.org/10.1080/0266476042000214501>.
- [9] D.E. Giles, *Superstardom in the us popular music industry revisited*, Econ. Lett. 92 (2006), pp. 68–74. <https://doi.org/10.1016/j.econlet.2006.01.022>.
- [10] S. Greven and T. Kneib, *On the behaviour of marginal and conditional AIC in linear mixed models*, Biometrika 97 (2010), pp. 773–789. <https://doi.org/10.1093/biomet/asq042>.
- [11] M. Hunger, A. Dring, and R. Holle, *Longitudinal beta regression models for analyzing health-related quality of life scores over time*, BMC Med. Res. Methodol. 12 (2012).



- [12] S.K. Kinney and D.B. Dunson, *Fixed and random effects selection in linear and logistic models*, *Biometrics* 63 (2007), pp. 690–698.
- [13] J. Lee and J.-S. Lee, *Music popularity: metrics, characteristics, and audio-based prediction*, *IEEE Trans. Multimedia* 20 (Nov 2018), pp. 3173–3182. <http://dx.doi.org/10.1109/TMM.2018.2820903>.
- [14] A. Lerch, *The Relation Between Music Technology and Music Industry*, Springer, Berlin Heidelberg, 2018, pp. 899–909. [http://dx.doi.org/10.1007/978-3-662-55004-5\\_44](http://dx.doi.org/10.1007/978-3-662-55004-5_44).
- [15] G. Lovison, M. Sciandra, A. Tomasello, and S. Calvo, *Modeling posidonia oceanica growth data: from linear to generalized linear mixed models*, *Environmetrics* 22 (2011), pp. 370–382. Available at <https://onlinelibrary.wiley.com/doi/abs/10.1002/env.1063>.
- [16] K. Middlebrook and K. Sheik, *Song hit prediction: predicting billboard hits using spotify data*, 2019.
- [17] M. Nasreldin, *Song popularity predictor*. Available at <https://towardsdatascience.com/song-popularity-predictor-1ef69735e380>, 2018.
- [18] Y. Ni, R. Santos-rodriguez, M. Mcvicar, and T.D. Bie, *Hit song science once again a science?* 2015.
- [19] R. Nijkamp, *Prediction of product success: explaining song popularity by audio features from spotify data*, July 2018. Available at <http://essay.utwente.nl/75422/>.
- [20] F. Pachet, *Musical metadata and knowledge management*, in *Encyclopedia of Knowledge Management*, 2nd ed., D. Schwartz, and D. Te'eni, eds., IGI Global, Sony CSL, Paris, France, 2011, pp. 1192–1199.
- [21] M. Poggini, *Liga. La biogra a*. BUR Biblioteca Univ. Rizzoli, ITALIA, 2010. ISBN 10: 8817040053.
- [22] M. Sciandra and A. Plaia, *A graphical model selection tool for mixed models*, *Comm. Stat. Simul. Comput.* 47 (2018), pp. 2624–2638. <https://doi.org/10.1080/03610918.2017.1353617>.
- [23] J.A. Sloboda, *Music in everyday life, the role of emotions*, in *Handbook of Music and Emotion: Theory, Research, Applications*, P. N. Juslin and J. Sloboda, eds., Oxford University Press, Oxford, 2011, pp. 1–37.
- [24] SpotifyWebAPI, *Spotify for developers*, 2019. Available at <https://open.spotify.com/>.
- [25] F. Vaida and S. Blanchard, *Conditional Akaike information for mixed-effects models*, *Biometrika* 92 (2005), pp. 351–370. <https://doi.org/10.1093/biomet/92.2.351>.