# Feature Engineering for Genre Characterization in Brazilian Music

Bruna Wundervald[1][0000−0001−8163−220X]

Maynooth University, `bruna.wundervald@mu.ie`

**Abstract.** Many factors are involved in the definition of music genres, making it an active area of research. This work focuses on verifying the connection between harmonic information and genre specification in Brazilian music, through the evaluation of feature importance in machine learning models. We construct four different sets of manually engineered harmonic features and assess how they relate to the accuracy of the models, as well as explore the mistakes made by the model in each genre. We identified the most relevant features to be the harmonic ones, followed by external features such as popularity and the proportions of the most common chord transitions in each song.

**Keywords:** chord features · feature importance · genre characterization.

## 1   Introduction

Genre is an important form of classifying songs, as they facilitate the search for music, and users even prefer to use genre instead of other metrics when looking for new music [9]. However, many factors are involved in the configuration of a music genre, such as style, historical context, and harmonic structures [3], making the definition of each genre unclear. Inconsistencies and blurriness in the definition of musical genres pose an important problem in various aspects of music studies and is an active area of research in MIR. For such reasons, the focus of this work is towards verifying the connection between harmonic information and genre specification in Brazilian music through the evaluation of feature importance in machine learning models. In addition, as [4] and [6] observed, mid-level music features such as chords configure a rich resource of information regarding genres. The chords sequence of a song fully describes its harmonic progression and it represents a meaningful part of the total music structure. With that, in this work, we also focus on the use of symbolic chords data and in manually extracting harmonically related features for genre classification, representing the chords structures in different and meaningful forms.

Related work has been done all of the usual representations of music data, for example, [14], [18], [1] and [17], which focused in music genre classification using audio extracted features. As for text data related to music, [11] presents a discussion about the characterization of genres through song lyrics. In [5], the authors introduced a vector based representation for chords sequences, bringing

light to an effective way to extract information about symbolic chords data. A similar problem to ours was studied in [15], which focused on harmonic features for genre classification.

## 2 Definitions

### 2.1 Data

The data was extracted from the Cifraclub website (https://www.cifraclub.com.br/), an online collaborative page of music-sharing, via the `chorrrds` [19] package for `R` [16]. Though the use of user-inputed (or crowd-sourced) music chords is not very common in MIR due to the possible inconsistencies, recent literature [12, 7] has been showing its value to the community. In total, 8 music genres were used: Reggae, Pop, Forró, Bossa Nova, Sertanejo, MPB, Rock and Samba, all good representatives of the Brazilian music, and from these genres, 106 different artists were available in the online platform, for which the chords and keys for 8339 different songs were collected. Complementary features about the release year and popularity were obtained with the aid of the well-known Spotify *API*.

### 2.2 Manually Extracted Features

In this work, we emphasized on obtaining various interpretable summary features from the chords, to make use of more information than only the symbolic form of the chords. The engineered features were separated into four thematic groups, organized as the **First set, triads and simple tetrads:** percentage of suspended chords (e.g. Gsus), of chords with the seventh (e.g. C7), of minor chords with the seventh (e.g. Em7, C#m7), of minor (e.g. Em, C#m), of diminished (e.g. Bᵒ), and of augmented (e.g. Baug) chords. **Second set, dissonant Tetrads:** percentage of chords with the fourth (e.g. D4), the sixth (e.g. E6), the ninth (e.g. G9), with the major seventh (e.g. F7+, Am7+), with a diminished fifth (e.g. C5- or C5b) and with an augmented fifth (e.g. C5+ ou C5#). **Third set, main chord transitions:** percentage of the first, second, and third most common chord transitions in the song. **Fourth set, miscellany:** popularity, total of non-distinct chords, year of album release, indicator of the key of the song being the same as the most common chord, percentage of chords with varying bass (e.g. C/E, C/G, C/Bb), mean distance of the root note to 'C' in the circle of fifths, mean distance of the root note to 'C' in semitones, absolute number of the most common chord.

### 2.3 Machine Learning Algorithm

We used the popular Random Forest [2] model, which is mainly characterized by being a tree ensemble that only allows a random subset $m$ of the features to be the candidates for a split, helping to create uncorrelated trees. This bagged ensemble can be written as $\hat{f}(\mathbf{x}) = \sum_{n=1}^{N_{tree}} \frac{1}{N_{tree}} \hat{f}_n(\mathbf{x})$, where $\hat{f}_n$ corresponds to the $n$-th tree.

# 3   Results

**Table 1.** Goodness of fit for the four models: overall accuracy with lower and upper bounds and Kappa statistic with the respective p-value.

| Model | Accuracy | L.B. | U.B. | Kappa | P-Value |
|---|---|---|---|---|---|
| Model 1 | 0.53 | 0.51 | 0.55 | 0.37 | < 0.0001 |
| Model 2 | 0.57 | 0.54 | 0.59 | 0.42 | < 0.0001 |
| Model 3 | 0.59 | 0.56 | 0.60 | 0.44 | < 0.0001 |
| Model 4 | 0.62 | 0.60 | 0.64 | 0.49 | < 0.0001 |

Our target variable here is the music genres, and the predictors are the engineered features. There is extensive literature in genre classification and we do not intend to claim that this is a better model than the others, as our primary goal is to observe how the features relate to the accuracy rather than obtaining the best accuracy possible. Four models were fitted in a nested fashion, with each new model being added with one of the features sets described before. Table 1 shows that, for all different models, there is evidence of their accuracy being significantly higher the non-information classification rate. The addition of feature sets progressively increases the accuracy of the models, evidencing that the 4 sets of features are informative to predict the genres. The increase is seemingly uniform: to each new set of variables added, the increase is about 3%.

**Table 2.** Confusion matrix for the model with all the features.

| | Bossa Nova | Forró | MPB | Pop | Reggae | Rock | Samba | Sertanejo |
|---|---|---|---|---|---|---|---|---|
| Bossa Nova | **0.28** | 0.00 | 0.40 | 0.00 | 0.00 | 0.05 | 0.16 | 0.12 |
| Forró | 0.00 | **0.00** | 0.12 | 0.00 | 0.00 | 0.12 | 0.10 | 0.65 |
| MPB | 0.01 | 0.00 | **0.59** | 0.00 | 0.00 | 0.11 | 0.13 | 0.15 |
| Pop | 0.00 | 0.00 | 0.13 | **0.00** | 0.00 | 0.28 | 0.15 | 0.44 |
| Reggae | 0.00 | 0.00 | 0.25 | 0.00 | **0.08** | 0.46 | 0.08 | 0.12 |
| Rock | 0.00 | 0.00 | 0.16 | 0.00 | 0.00 | **0.43** | 0.05 | 0.35 |
| Samba | 0.01 | 0.00 | 0.20 | 0.00 | 0.00 | 0.03 | **0.66** | 0.10 |
| Sertanejo | 0.00 | 0.00 | 0.02 | 0.00 | 0.00 | 0.07 | 0.02 | **0.89** |

Figure 1 shows that the first set of features is the most informative one, meaning that with the basic chords information we can already obtain good results in terms of informing the model about the genres. The external variables, such as the year and popularity, got a high rank in the plot, showing how the Spotify features are also pertinent. The position of the transitions and distances variables strengthen the idea of harmonic characteristics being important to discriminate
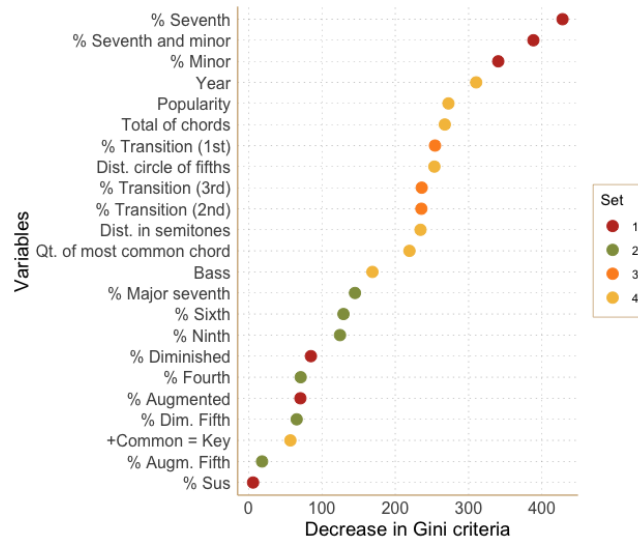
**Fig. 1.** Importance plot for the fourth model with all the considered features. The top part of the plot is dominated by harmonic features.

music genre. From Table 2, we can see that there is considerable confusion between MPB and Bossa Nova, highlighting their known harmonic similarities. The same happens to Forró, Sertanejo and Pop, which are music genres with a similar origin and, in general, more elementary harmonic structures.

## 4   Conclusions

With our results, we conclude that manually engineered harmonic features can be useful to to characterize Brazilian music genres. More than just predicting music genres, which does not have a consensual utility in the literature, we are interested in inferring which harmonic features are informative for the definition of genres. In our case, the most discriminative features are the percentage of chords with the seventh note, of minor chords with the seventh note, of minor chords, the year of release of the songs, the popularity and the behavior of the most common chord transitions. Apart from that, though our work was limited to one geographic region, we believe that our insights can be extended to other types of music that influenced or were influenced by the genres considered here, such as Jazz, Pop, and Rock music.

The next steps of this work include specially the engineering of the new variables and applying different algorithms, such as deep learning models [8] and naive Bayes models [10], as in [1] and [13], though they might have less interpretable results. In a different sense, we would also like to explore more the use of crowd-sourced data and the relationship between song popularity and the precision of this type of data.

## References

1. Bahuleyan, H.: Music genre classification using machine learning techniques. CoRR **abs/1804.01149** (2018), http://arxiv.org/abs/1804.01149
2. Breiman, L.: Random forests. Machine Learning (2001). https://doi.org/10.1023/A:1010933404324
3. Caldas, W.: Iniciação à Música Popular Brasileira, vol. 1 (2010)
4. Cheng, H.T., Yang, Y.H., Lin, Y.C., Liao, I.B., Chen, H.H.: Automatic chord recognition for music classification and retrieval. In: 2008 IEEE International Conference on Multimedia and Expo. pp. 1505–1508. IEEE (2008)
5. Chuan, C.H., Agres, K., Herremans, D.: From context to concept: exploring semantic relationships in music with word2vec. Neural Computing and Applications **32**(4), 1023–1036 (2020)
6. Corrêa, D.C., Rodrigues, F.A.: A survey on symbolic data-based music genre classification. Expert Systems with Applications **60**, 190–210 (2016)
7. Koops, H.V., de Haas, W.B., Bransen, J., Volk, A.: Automatic chord label personalization through deep learning of shared harmonic interval profiles. Neural Computing and Applications **32**(4), 929–939 (2020)
8. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. nature **521**(7553), 436–444 (2015)
9. Lee, J.H., Downie, J.S.: Survey of music information needs, uses, and seeking behaviours: preliminary findings. In: ISMIR. vol. 2004, p. 5th. Citeseer (2004)
10. Murphy, K.P., et al.: Naive bayes classifiers. University of British Columbia **18**, 60 (2006)
11. Neuman, Y., Perlovsky, L., Cohen, Y., Livshits, D.: The personality of music genres. Psychology of Music (2016). https://doi.org/10.1177/0305735615608526
12. Odekerken, D., Koops, H.V., Volk, A.: Decibel: Improving audio chord estimation for popular music by alignment and integration of crowd-sourced symbolic representations. arXiv preprint arXiv:2002.09748 (2020)
13. Oramas, S., Nieto, O., Barbieri, F., Serra, X.: Multi-label music genre classification from audio, text, and images using deep features. CoRR **abs/1707.04916** (2017), http://arxiv.org/abs/1707.04916
14. Pampalk, E., Flexer, A., Widmer, G.: Improvements of Audio-Based music similarity and genre classification. In: ISMIR (2005). https://doi.org/10.1007/s10115-013-0641-y
15. Pérez-Sancho, C., Rizo, D., Iesta, J.M., De León, P.J., Kersten, S., Ramirez, R.: Genre classification of music by tonal harmony. Intelligent Data Analysis (2010). https://doi.org/10.3233/IDA-2010-0437
16. R Core Team: R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria (2018), https://www.R-project.org/
17. Scaringella, N., Zoia, G., Mlynek, D.: Automatic genre classification of music content. IEEE Signal Processing Magazine (2006). https://doi.org/10.1109/MSP.2006.1598089
18. Tzanetakis, G., Cook, P.: Musical genre classification of audio signals. IEEE Transactions on Speech and Audio Processing (2002). https://doi.org/10.1109/TSA.2002.800560
19. Wundervald, B.: The chorrrds package for extraction of music chords data in r (2018), https://github.com/r-music/chorrrds