

Explorando o Populismo nos Discursos de Jair Bolsonaro: Um Estudo Estatístico

Isabela do Canto – Departamento de Estatística - UFPE

Isabela.canto@gmail.com

<https://orcid.org/0009-0002-0720-3014>

Maria do Carmo Soares de Lima – Departamento de Estatística - UFPE

maria@de.ufpe.br

<https://orcid.org/0000-0002-5480-3103>

Resumo - Este estudo busca fornecer uma análise inicial utilizando um corpus derivado dos discursos proferidos durante o primeiro ano do governo do ex-presidente Jair Bolsonaro. Antes de avaliar o desempenho de vários algoritmos de aprendizagem automática amplamente utilizados, foram aplicadas técnicas básicas de pré-processamento dos dados. Os resultados mostraram que o algoritmo Support Vector Classification alcançou o melhor desempenho, com uma taxa de precisão de 86%. Além disso, o estudo incluiu uma análise baseada em bigramas para distinguir e comparar os padrões linguísticos entre os discursos classificados como populistas e aqueles considerados não populistas. Esses resultados iniciais oferecem insights valiosos sobre as estratégias retóricas empregadas no início da presidência de Bolsonaro.

Keywords: Jair Bolsonaro; Machine Learning; Support Vector Classification algorithm

Exploring Populism in Jair Bolsonaro's Speeches: A Statistical Study

Abstract – This study seeks to provide an initial analysis using a corpus derived from the speeches delivered during the first year of former President Jair Bolsonaro's administration. Before assessing the performance of several widely-used machine learning algorithms, basic pre-processing techniques were applied to the data. The findings revealed that the Support Vector Classification algorithm achieved the highest performance, with an accuracy rate of 86%. Furthermore, the study included an analysis based on bigrams to distinguish and compare the linguistic patterns between speeches classified as populist and those considered non-populist. These initial results offer valuable insights into the rhetorical strategies employed during Bolsonaro's early presidency.

Keywords: Jair Bolsonaro; Machine Learning; Support Vector Classification algorithm.

Data da Submissão: 16/08/2024

-

Data de aceitação: 10/01/2025

Acknowledgments: Special thanks to Davi Moreira and Caio Vinicio Malaquias do Vale for making the data available.

Os direitos autorais são atribuídos às pessoas autoras do artigo.

Este artigo está licenciado sob forma de uma licença Creative Commons Atribuição-Não Comercial-Sem Derivações 4.0 Internacional (CC BY-NC-ND 4.0).
<https://creativecommons.org/licenses/by-nc-nd/4.0/>



1. Introduction

In recent years, interest in the rise of populism has surged, leading to various branches of study. The most relevant of these is the relationship between populism and the legitimacy crisis faced by several industrial democracies, which were previously considered to be consolidated (Moffit, 2016). The loss of support for traditional institutions such as political parties and the media presents, in this context, a window of opportunity for the rise of a new radical right (Liñán, 2018). The study of political actors that align with this radical right is, therefore, crucial for a better understanding of this phenomenon. For this reason, the central objective of this article is to conduct a statistical analysis of the speeches from the first year of former President Jair Bolsonaro's term, who is classified by many as a populist leader (Ricci, Izumi, and Moreira, 2021).

This article aims to explore the discursive strategies employed by former Brazilian President Jair Bolsonaro, focusing on identifying populist elements within his speeches. Through the application of machine learning techniques, the study seeks to offer a dual contribution: first, to enhance understanding of Bolsonaro's rhetoric in the context of populism, and second, to advance the methodological discussion on the use of classification algorithms in political discourse analysis. The gap in the literature lies in the limited focus on the intersection of political rhetoric and advanced technological tools, and this work aims to fill that by using a diverse set of machine learning models.

However, the goal is not to simply classify the discourse as populist or non-populist, but rather to analyze the effectiveness of different techniques for this purpose using a dataset that has already been manually classified, as well as to conduct a deeper analysis of the characteristics present in this type of discourse. This methodological discussion is crucial for developing tools that assist researchers in studying populism, which is often conducted holistically across various stages. This does not delegitimize the performance of the holistic approach but instead offers alternative analytical perspectives that contribute to a broader understanding through the application of new tools. This objective will be pursued through the application of statistical models to a corpus of 140 sentences from the year 2019, which have been previously classified as “populista” (populist) or “não-populista” (non-populist) by experts in the field.

2. Conceptual Framework

Since classification is an abstract activity, there is a need for clarity in its definition to ensure it is as transparent and replicable as possible. This justifies using the definition of populism as the construction of polarization in society, seen as an opposition between “povo” (people) and “elite” (Hawkins, 2017), which is justified by social and economic cleavages that vary according to the regional context. The “elites” also vary depending on the ideological spectrum of the populist, ranging from political, economic, legal, and supranational elites to the media. As pointed out by Engesser (2016), left-wing populists tend to focus their discourse on economic elites, while right-wing populists often target the media. In the case of Brazil, a prominent elite for the far-right is the Supreme Court (Supremo Tribunal Federal or STF), which is considered a limiter of freedom (limitador da liberdade).

Bolsonaro’s discourse aligns with this populist narrative, as previous studies have demonstrated (Ricci, Izumi, and Moreira, 2021). The theoretical framework in this article connects this populist paradigm to Bolsonaro's speeches and uses advanced machine learning techniques to (1) classify these elements and (2) analyze the pattern of the terms used, thereby gaining a deeper understanding of the strategic patterns present in the former president's discourse. The focus is on addressing how modern technology can be leveraged to interpret traditional political phenomena.

Understanding the need for political discourse analysis in the current democratic crisis is essential, given that it is a phenomenon that relies on “performance” to construct societal divisions by emphasizing social and/or economic crises and anti-establishment rhetoric that delegitimizes traditional institutions (Moffitt, 2015). Since this discursive construction depends on the political, economic, and social context of the country, the elements mobilized vary. This explains the coexistence of six major theoretical families: Political Leadership, Political Culture, Political Institutions, Political Economy, Social Structure and Political Coalitions, and International Factors (Lust & Waldner, 2015a). The analysis of presidential discourse is based on the premise of the predominance of the Political Leadership theory in explaining the phenomenon, which, while necessary, is not sufficient. A thorough analysis of the elements mobilized within the speeches reveals the causal complexity by presenting variables from the Economic Theory, Social Structure, and Political Institutions theories.

The relevance of Jair Bolsonaro's discourse is justified by Brazil's personalized presidential nature, but especially by the new mobilizing behavior adopted by populist/extremist actors, characterized by an “audience-based populism” that relies on discourse directed at the “povo” (people) to legitimize authoritarian actions. Many studies analyze the use of social media among populists, such as Engesser (2017). The authors consider that:

“Five key elements of populism are derived from the literature: emphasizing the sovereignty of the people, advocating for the people, attacking the elite, ostracizing others, and invoking the ‘heartland’. A qualitative text analysis reveals that populism manifests itself in a fragmented form on social media”. (Engesser, 2017, p 8).

Although this article does not analyze social media, the same elements are used in official speeches, as will be presented later.

3. Methodology

This study employed a mixed approach to political discourse analysis, combining text preprocessing techniques with machine learning algorithms to classify populist and non-populist speeches. The dataset used was derived from Ricci (2021), consisting of speech sentences that met the definition of populism as a dichotomy between the people and the elite (Hawkins, 2017). The initial classification was manually conducted by experts, using a dictionary of keywords related to the concepts of "people" and "elite."

Text preprocessing was a crucial step to ensure the effectiveness of the classification models. Techniques such as lemmatization and the removal of stop-words, including pronouns and articles, were applied to standardize and clean the text. Furthermore, tokenization and grammatical categorization methods were employed to refine the linguistic analysis.

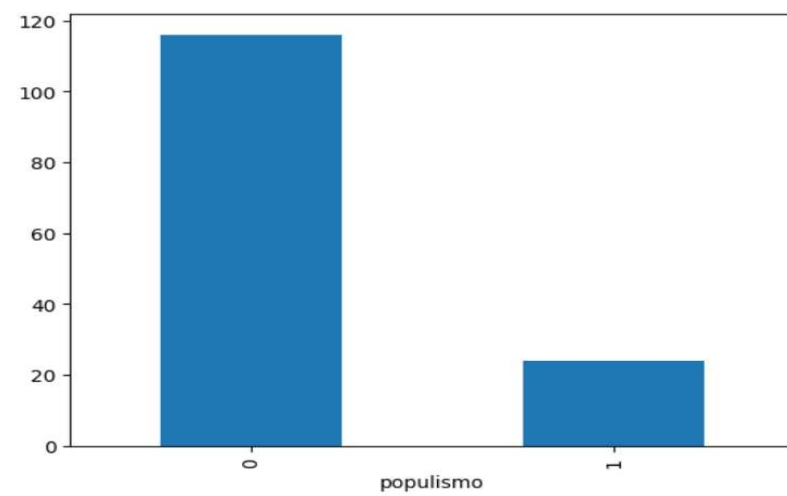
Once the text was processed, several machine learning algorithms were tested. The Support Vector Classification (SVC) model demonstrated the best performance in distinguishing between populist and non-populist speeches. However, to ensure robust results, other algorithms such as Naive Bayes, KNN, and Neural Networks were also employed, allowing for a comparative analysis of different techniques.

Performance metrics evaluated included precision, recall, F1-score, and accuracy, aiming to assess the models' ability to differentiate between the two speech categories. The dataset was split into 85% for training and 15% for testing, providing a robust sample to validate the findings. Finally, a comparative analysis was conducted, considering bigrams and TF-IDF metrics to identify the most frequent words and expressions in both types of speeches.

3.1 Data collection

The database used in this article is from Ricci (2021), consisting of speech sentences that, following the definition of populism as an opposition between the people and the elite (Hawkins, 2017), met the criterion of containing words from a dictionary that expresses elements referring to the idea of "elite" and elements referring to the idea of "people." In addition to the dictionary application, the resulting sentences were manually classified as populist or non-populist by experts. This process resulted in a database of potentially populist sentences, which were classified with both positive and negative outcomes that were addressed in this study.

Figure 1: Descriptive statistics for the data.



From the data set used, 140 speeches are from former president Jair Bolsonaro, of which, using the classification of populism (1) and non-populism (0) adopted by Ricci, Izumi and Moreira (2021), 116 of them were identified as non-populist, which corresponds to approximately of 83%, as showed in figure 1.

Table 1 shows the number of documents in which each of the words related to “povo” appears, as well as their respective representations (in terms of percentage) in speeches considered populist and non-populist.

3.2 Text preprocessing

In order, to perform an analysis of the collected speeches, the data was pre-processed using machine learning (ML) techniques. This step is important for the analysis, as pointed out by some works such as Joshi (2020) and Torfi et al. (2020). Therefore, the following basic techniques were used:

- Using lemmatization techniques, the words returned to their root form;
- Through tokenization, possessive pronouns and articles were removed, that is, we depersonalized the subject to make the text “cleaner”. Furthermore, stop-words were removed, such as auxiliary verbs, determiners and conjunctions;
- The grammatical classes (or Part of Speech - POS) were also taken into account in the pre-processing of the corpus in question. To do so, the Spacy package and its functions were used. Table 1: Checking the number (in percentage) of keywords related to the word “povo” considered populist and non-populist present in the documents.

Table 1: Checking the number (in percentage) of keywords related to the word “povo” considered populist and non-populist present in the documents.

| Keyword | Class | 1 | n | Keyword | Class | 1 | n |
|------------------|-------|-------|----|-----------------|--------|-------|----|
| “amiga” | 67.00 | 33.00 | 3 | “nordestina” | 0.00 | 0.00 | 0 |
| “amigas” | 0.00 | 0.00 | 0 | “nordestinas” | 0.00 | 0.00 | 0 |
| “amigo” | 2.41 | 97.59 | 83 | “nordestino” | 77.78 | 22.22 | 9 |
| “amigos” | 0.00 | 0.00 | 0 | “nordestinos” | 0.00 | 0.00 | 0 |
| “brasileiros” | 0.00 | 0.00 | 0 | “Norte” | 42.86 | 57.14 | 7 |
| “brasileiras” | 0.00 | 0.00 | 0 | “pobre” | 65.22 | 34.78 | 23 |
| “classe média” | 0.00 | 0.00 | 0 | “pobres” | 0.00 | 0.00 | 0 |
| “classe popular” | 0.00 | 0.00 | 0 | “pobreza” | 75.00 | 25.00 | 4 |
| “cristã” | 80.00 | 20.00 | 40 | “povo” | 100.00 | 0.00 | 1 |
| “cristão” | 78.95 | 21.05 | 38 | “trabalhador” | 81.82 | 18.18 | 11 |
| “cristãos” | 0.00 | 0.00 | 0 | “trabalhadora” | 50.00 | 50.00 | 2 |
| “cristãs” | 0.00 | 0.00 | 0 | “trabalhadoras” | 0.00 | 0.00 | 0 |
| “família” | 77.81 | 23.19 | 69 | “trabalhadores” | 0.00 | 0.00 | 0 |
| “Nordeste” | 0.00 | 0.00 | 0 | | | | |

Source: Authored by the researchers

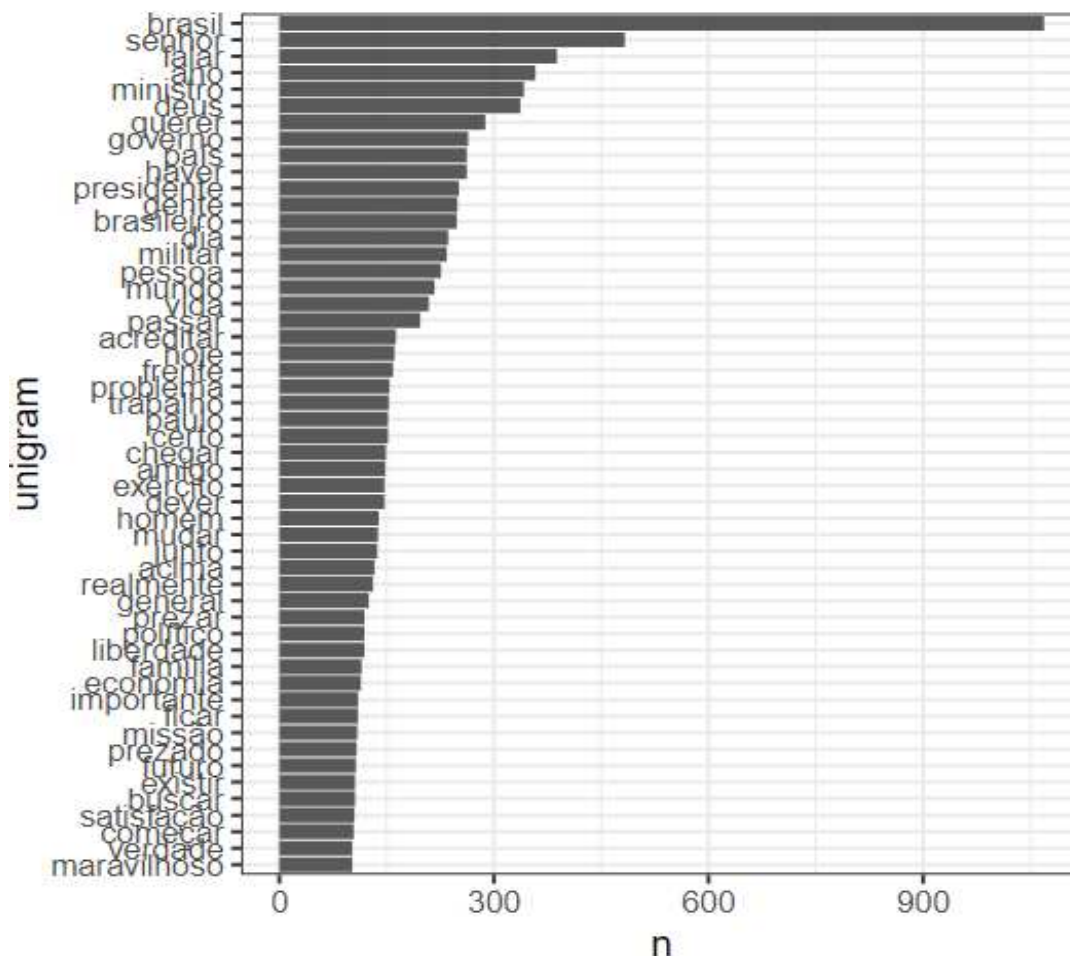
After carrying out the pre-processing techniques described above, the new corpus contains 64,922 words, of which 6,495 were found in speeches considered non-populist. To provide a simplified portrait, figure 2 shows the words that have a frequency greater than 100 in speeches considered non-populist.

To understand the frequency of words in speeches considered populist, we illustrate in the figure 3 those that appeared more than 40 times in the speeches.

When comparing figures 2 and 3, it is possible to notice the presence of the word ‘Brasil’ as the one that appears most in both speeches considered non-populist and populist, more frequently in the first.

However, it is worth highlighting that the word ‘brasileiro’ appears in both types of speeches, as can be seen in the figures above, appearing 248 times in speeches considered non-populist and 99 times in speeches considered populist.

Figure 2: Word frequency in speeches considered non-populists.



3.3 Populism dictionary

The classification considered here was the same as that adopted by Ricci, Izumi and Moreira (2021), which was done manually.

It is important to note that the construction of this dictionary was specifically designed to suit the context of populism in Brazil, incorporating elements that are mobilized within these political discourses. For this reason, it is crucial to emphasize the need for adapting certain terms for different applications, as the elements that characterize populism heavily depend on political divisions, social structure, and the specific challenges faced by each society.

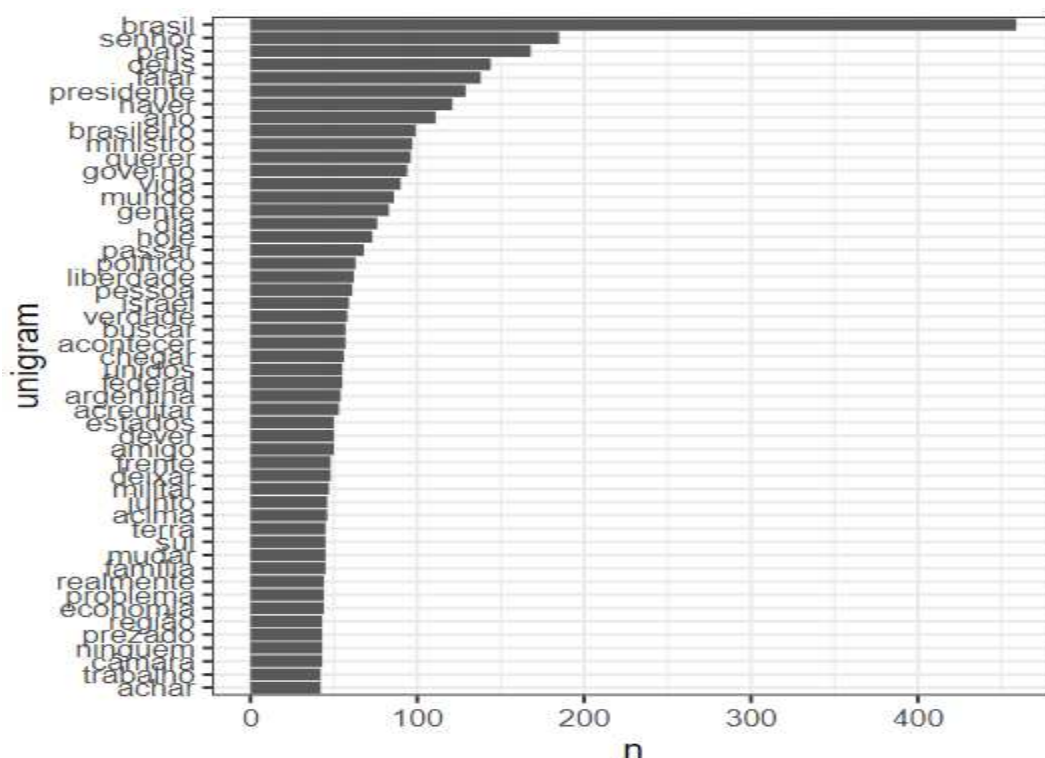
The case of Brazil, in particular, requires an analysis that captures the patterns of religiosity, conservatism, and nationalism present among different groups. Therefore, terms such as “Deus” (“God”), “Israel,” “militar” (“military”), “família” (“family”), “liberdade” (“freedom”), and “ministro” (“minister”) are included. These terms reflect the current debates that polarize political groups in their rhetoric.

The first two terms are particularly relevant to the evangelical group; the third signifies the authoritarian tendency perceived as necessary to address societal problems;

the fourth reflects a defense of conservatism in opposition to the liberal ideas advocated by the center-left; the fifth represents the rhetoric that defends the freedom to express this authoritarian tendency; and the last one signifies the anti-establishment element, often through attacks against the Supreme Court (STF), which further contributes to the authoritarian tendency.

By comparing figures 2 and 3, we observe that many words appear in both populist and non-populist statements. This is the case with the term “Brasil” (“Brazil”). Such overlaps highlight a limitation of simpler techniques for the classification of discourses.

Figure 3: Word frequency in speeches considered populists,



3.4 Metrics and techniques for resource extraction

Here, we consider one word related to “people” and one word related to “elite” to study the TF-IDF. The choice was made randomly, considering a minimum number of 10 bigrams, as some words characterized as relating to people and elite could be absent in the corpus. For both cases, two attempts were made until finding (in the second) the minimum number. The results are presented in Table 2 and 3.

The following technique was used:

- Bigrams were considered by taking the second word as the term considered and we treated the bigram as just one word;
- The TF-IDF statistic was calculated.

Tables 2 e 3 show, in descending order, the bigrams that are most relevant in the corpus under study, according to the TF-IDF metric.

It is possible to note that the bigrams that have the highest TF-IDF, 2.8903, are those related to the “elite” (elite), totaling 14 bigrams. In addition, all of these 14 bigrams have the word “empresário” (entrepreneur) as the second word of the bigram. On the other hand, the bigrams with the highest TF-IDF of the words related to “povo” (people) have “nordestino” (northeastern) as the second word.

The Northeast is a region of most importance for Bolsonaro, both due to the number of voters—being the second largest region, with 27.7% of the electorate—and the difficulty in securing votes, as it is a region that traditionally supports the left-wing Workers’ Party (PT) for various reasons. These reasons range from the socio-economic structure and the benefits the population receives from the public policies of income redistribution and poverty reduction implemented by the party to the predominant religion. Former President Bolsonaro has stronger support among evangelical voters, while the region is predominantly Catholic (72%) and has the lowest percentage of evangelicals (16%) in the country (Nicolau, 2018).

As a strategic region for electoral gains, it cannot be overlooked in discourse, justifying the expressions adopted by Bolsonaro that aim to connect with the public in this region, as we can see in Table 2.

Another significant topic in the rhetoric of this former president is his relationship with business leaders, which explains the results presented in Table 3. As a far-right politician, his connection with the country’s economic elite is particularly noteworthy. By advocating neoliberal policies, he garnered support from prominent figures such as Luciano Hang (Lojas Havan), Flávio Rocha (Grupo Guararapes/Lojas Riachuelo), João Apolinário (Polishop), among others.

Table 2: BoW representation with TF-IDF - Keywords relating to “povo”.

| | Bigram | TF-IDF |
|---|--|--------|
| 1 | “dória nordestino”, “filho nordestino” “gostar nordestino”, “sertão nordestino” | 2.1972 |
| 2 | “semiárido nordestino” | 1.5040 |
| 3 | “roubar nordestino” | 1.0986 |
| 4 | “irmão nordestino” | 0.7520 |
| 5 | “cidade nordestino” “falar nordestino” | 0.7324 |
| 6 | “abraço nordestino”, “beijo nordestino” “ceará nordestino”, “hoje nordestino” | 0.5493 |
| 7 | “irmão nordestino” | 0.5013 |

Source: Authored by the researchers.

Table 3: BoW representation with TF-IDF - Keywords relating to “elite”.

| | Bigram | TF-IDF |
|---|---|--------|
| 1 | “cientista empresário”, “presente empresário” “havan empresário”, “americano empresário” “procurar empresário”, “ano empresário” “reunir empresário”, “provisória empresário” “obrigatoriedade empresário”, “país empresário” “amigo empresário”, “doação empresário” “empreendedores empresário”, lembrar empresário | 2.8903 |
| 2 | “conversar empresário”, “nome empresário” “agradecemos empresário”, “confiança empresário” | 1.4451 |
| 3 | “mão empresário”, “senhores empresário” “brasil empresário”, “veja empresário” | 0.9634 |
| 4 | “senhor empresário” | 0.7324 |

Source: Authored by the researchers.

3.5 Experimental setup

It's very important.

3.5.1 Dataset division

The data set under study was divided into 85% for training and 15% for testing, which corresponds to 119 and 21, respectively. These values for test and control were chosen due to tests performed. This scenario was the one that generated the best results.

3.5.2 Performance metrics

In order, to evaluate the models studied here, some commonly used metrics were considered, namely: precision, recall, f1-score, and accuracy. The first of them is the fraction of relevant instances among the retrieved instances. The recall metric describes the fraction of true positives and all positives, that is, the higher the value, the better the performance of the algorithm is to capture true positive elements.

The third metric used here, called f1-score, is a balanced combination of the first two, and is therefore defined as the harmonic mean between the first two metrics. Therefore, a good model, considering the f1-score metric, will be able to simultaneously correct its predictions and retrieve the examples of the class of interest.

Finally, the last metric considered here, accuracy, can evaluate, among all the classifications, how many the model was able to classify correctly.

4. RESULTS

In this section, we consider the bigram-based study again. As there are many that can be formed from the corpus under study, which has 140 observations, we only take the bigrams that appear 50 or more times in the speeches.

The results are presented in Tables 4 and 5. For example, according to the results seen in the first table, it is possible to identify the two most cited bigrams in the documents: “forças armadas” and “paulo guedes”, appearing in 55% and 50.7% of the documents, respectively.

The results obtained seem reasonable in the context of the government of Jair Bolsonaro, for many reasons. The Bolsonaro administration frequently employed militaristic symbols and rhetoric, endorsing the view of this institution as a representative of nationalism, conservatism, order, and discipline. These characteristics were used to legitimize decisions to allocate military personnel to civilian positions, including many ministerial roles and public policy-making positions. This conferred fundamental relevance to the Armed Forces for understanding the government. The name ‘Paulo Guedes’ can also be justified by the fact that he was the Minister of Economy in the government of the president in question, from 2019 to 2023 and was a name used also to legitimize the economic decisions and results. Paulo Guedes gave several interviews during Bolsonaro’s government due to the pandemic context and was a figure present throughout the government.

Now, analyzing the second table, it is possible to see that the two most cited bigrams are ‘states united’ and again ‘paulo guedes’. Again, the results match the context. Note that during Bolsonaro’s government, there were a strong attempt to strengthen diplomatic relations on the part of Brazil, mainly due to the political and economic alignment between Trump and Bolsonaro. The brazilian president used the American president’s stance to legitimize his own nationalist, conservative, and climate-skeptical positions.

It is also possible to note that of the seven bigrams presented in Tables 4 and 5, five of them are present in both. However, the frequency with which they appear in speeches considered populist is greater than in speeches considered non-populist in all five cases.

In order, to better visualize the relationships between the bigrams that appear more frequently in the discourses under study, Figures 5 and 6 present a heat maps of the first 20 bigrams most frequently in discourses considered non-populist and populist, respectively.

Among the most significant relationships, there is a strong connection between the use of the words “Deus”, “Brasil” and “acima” in both figures, which is directly linked to Bolsonaro’s government motto “Brasil acima de tudo, Deus acima de todos”. Additionally, the previously mentioned influence of the armed forces is evident in the relationship between “military” and “academy” “school” and “agulhas” (a reference to a military institution); between “forças” and “armadas”; and between “exército” and “brasileiro”. The then-Minister of Economy, Paulo Guedes, is also represented in the heat maps, as expected, given his relevance during the administration.

Finally, we see a connection between the words "democracia" and "liberdade". It is not expected to appear in the non-populist heat map, as these words were often used to justify speeches that disregarded rights and/or institutions.

Table 4: Statistics of documents considered populist.

| word1 | word2 | percent |
|----------|------------|------------|
| forças | armadas | 0.55000000 |
| paulo | guedes | 0.50714286 |
| exército | brasileiro | 0.37857143 |
| rio | janeiro | 0.37142857 |
| deus | acima | 0.36428571 |
| brasil | acima | 0.35714286 |
| senhor | senhora | 0.35714286 |

Source: Authored by the researchers.

Table 5: Statistics of documents considered not populist.

| word1 | word2 | percent |
|---------|-----------|------------|
| estados | unidos | 0.35714286 |
| paulo | guedes | 0.15000000 |
| brasil | acima | 0.12142857 |
| deus | acima | 0.12142857 |
| forças | armadas | 0.12142857 |
| senhor | senhora | 0.10714286 |
| terras | indígenas | 0.10714286 |

Source: Authored by the researchers.

As previously noted, the TF-IDF bigrams failed to capture the relationship between words and populism. This is related to the difficulty in capturing the meaning of words within sentences. Being an abstract concept and with language itself being inherently challenging to interpret, the same expression can be used in both populist and non-populist sentences, depending on the meaning of the sentence and the role that other words play together. While the use of these metrics significantly contributes to the analysis of the former president's speeches, agendas, and interests, they are insufficient on their own to explain the occurrence of populism.

4.1 ML models's performance¹²

4.1.1 Naive Bayes with CV

Table 6: Results of models using different algorithms.

| Algotithm | Class | Precision | Recall | f1-score | Accuracy |
|--------------------|-------|-----------|--------|----------|-------------|
| CV | 1 | 0.84 | 0.94 | 0.89 | 0.81 |
| | 0 | 0.50 | 0.25 | 0.33 | |
| SVM | 1 | 0.00 | 0.00 | 0.00 | 0.86 |
| | 0 | 0.86 | 1.00 | 0.92 | |
| SVC | 1 | 0.89 | 0.94 | 0.91 | 0.86 |
| | 0 | 0.67 | 0.50 | 0.57 | |
| KNN | 1 | 0.20 | 0.33 | 0.25 | 0.71 |
| | 0 | 0.88 | 0.78 | 0.82 | |
| RF | 1 | 0.81 | 1.00 | 0.89 | 0.81 |
| | 0 | 0.00 | 0.00 | 0.00 | |
| NN(1) ² | 1 | 0.00 | 0.00 | 0.00 | 0.86 |
| | 0 | 0.86 | 1.00 | 0.92 | |
| NN(2) ³ | 1 | 0.00 | 0.00 | 0.00 | 0.76 |
| | 0 | 0.84 | 0.89 | 0.86 | |

Source: Authored by the researchers

The first algorithm that we will apply to the corpus under study is the Naive Bayes algorithm, one of the simplest and most efficient classifiers in the classification of data. Besides that, in order to evaluate the generalization of the Naive Bayes model, the cross-validation (CV) technique was used here.

From Table 6, it can be seen that 17 of 21 examples belong to class 1, which represents almost 81% of the documents, indicating an imbalance in the data set. Thus, the predictor in question almost always predicts any sample as belonging to class 1.

Besides that, we can therefore state that the data set is unbalanced. Therefore, let's analyze the confusion matrix illustrated in the following figure. It is possible to note that the algorithm makes a better prediction for “zeros” than for “ones”, which was already expected since the dictionary used by Ricci, Izumi and Moreira (2021) has a high classification of “zeros”.

Additionally, the Figure 7(a) reinforces what was discussed in the confusion matrix above. Figure 9e(a) shows that, as the sample size grows, the learning capacity increases, although the ability to make a correct prediction (true - positive) fluctuates just after $n = 50$.

¹ Neural networks method using 10 nêutrons and only one layer

² Neural Networks Method Using Multiple Layers.

4.1.2 SVM model

The Support Vector Machine (SVM) is a supervised machine learning method that has as one of its advantages being robust over overfitting, especially in high dimensionality. Although this is not the case in question, the objective here is to understand how the method manages to classify the discourses in question.

It is possible to notice that, in Figure 9e(b), as more training examples are used, the learning curve remains constant. On the other hand, as more data is added, the cross-validation curve grows to somewhere around 70, where it drops slightly. This drop after a certain point may indicate an *overfitting* problem. However, to reach a more accurate conclusion, it would be necessary to increase the test set for better analysis.

Regarding the standard deviation of the metrics used, the figure shows that there is a relatively large shaded area around the cross-validation curve, which indicates high variance between the different iterations of model training. However, it is possible to notice that this variance decreases from somewhere around 70.

4.1.3 SVC model

One of the most common ways to apply SVM to solve data categorization problems is by using Support Vector Classification (SVC). This algorithm has the advantages of being able to provide high accuracy and being less susceptible to overfitting.

The confusion matrix for the algorithm in question is plotted in Figure 7(c) and shows something quite similar to what was found using the CV algorithm. However, note that there is a difference only in the value of the correct classification of the “ones”.

The ROC curve for this case is shown in Figure 9b(b). Note that, for the use of this algorithm, all curves remained above the diagonal curve, differing from the behavior of the corresponding figure for the case of the CV algorithm. Thus, the figure indicates that the SVC model has the ability to distinguish between classes with some degree of ability.

Analyzing Figure 9e©, what we have is that the training curve remaining constant at 1 suggests that the model is fitting the training data very well. On the other hand, the oscillation facts in the cross-validation curve and the increase in the shaded area of this same curve indicate that this may be the result of different behavior in the different partitions of the cross-validation data. In short, the characteristics mentioned above indicate an overfitting challenge.

4.1.4 KNN model

Considering the context of the k-Nearest Neighbors algorithm (denoted here as KNN), the main parameter is the number of nearest neighbors (k). Here, we initially select $k = 5$. The results can be seen in Table 6 and Figures 7(d) and 9b(c), which

represent, respectively, the confusion matrix and the ROC curves.

Firstly, it is possible to notice that the method is more accurate in non- populist speeches than in populist speeches, as expected and similar to what was found with the use of the two algorithms above. In the second figure, it is noted that the ROC curves for both classes coincide. The fact that they are the same indicates that the model's performance is consistent in terms of its ability to distinguish between both classes. There is no significant disparity in performance between class 1 rating versus class 0 rating.

Besides that, Figure 9e(d) shows that, as the sample size grows, the learning capacity increases, although the ability to make a correct prediction (true - positive) fluctuates.

It is also possible to note that the variability or uncertainty associated with the performance of this algorithm seems to represent a decrease in the learning curve as the size of the training sample increases. However, the oscillation of the cross-validation curve indicates that the performance of this model can vary significantly, therefore suggesting that the model is not generalizing well. These aspects often indicate a sign of overfitting.

4.1.5 RF model

The results of the metrics for this model are also found in Table 6. However, this value does not reflect the good performance of the algorithm, since the true positives, i.e., the numbers of examples "one" correctly classified as "one" is zero. This fact ends up influencing a high accuracy value, even without a good performance of the model. In short, this result was already expected since this type of algorithm does not deal well with class imbalances.

The confusion matrix for this case, which can be seen in the Figure, shows us that the model is extremely biased, presenting a poor performance for the minority class, that is, for the populist class.

This type of situation is expected when the model has a high bias or when there is an extreme imbalance between classes. It is believed that, in this case study, the problem lies in the high classification of discourses as being non- populist.

Figure 9e(e) indicates that the model can correctly predict the vast majority of training examples, which means that the model is close to its maximum learning capacity. However, as the cross-validation curve experiences a slight increase just after $n=70$, we can say that the model is still able to generalize better to a larger number of training data.

4.1.6 NN model

Another model tested here is the Neural Networks (NN). Here, we use the "binary cross entropy" loss function for binary classification, with "adam" optimization method. The results obtained are shown in Table 6.

It is possible to note that all metrics obtained for the class considered pop- ulist

were null, as can be confirmed in Figure 7(f), which represents the confusion matrix for this case. In other words, the model failed to get any of the examples of the class considered populist right. This can be justified by the fact that the set has only 24 speeches considered populist out of a total of 140 speeches, which causes a possible problem of class imbalance. This fact may lead the model to have difficulty learning to distinguish examples of the class considered populist. The classifier used for the production of the Figure 9e(f) was the “MLP- Classifier” of the Python scikit-learn library. The acronym “MLP” refers to “Multilayer Perceptron”, which is a type of artificial neural network with multiple layers of neurons. In this case, the hidden layer was considered to have a number of 10 neurons and only one layer. There was no justification for choosing this number.

In addition, for the production of the learning curve, the training of the model with different subsets of the training set (10%, 20%, ... up to 100% of the data). With each iteration, the model is trained with an increasing amount of data, and training and testing accuracy is recorded.

As can be seen in Figure 9e(g), initially there is a sharp drop in both curves. Then the two curves fluctuate a lot with the training score curve oscillating less than the test score curve. Also note that at about the value of the abscissa axis at 45, the test curve is slightly above the training curve.

The above facts can indicate several things. One of these is the problem of underfitting, which occurs when the model is too simple to capture the complexity of the data. The impact seen in the graph can perhaps be justified by the fact that the neural network has only one layer and 10 neurons.

In view of this, it was decided to treat this possible problem by increasing the number of layers and neurons per layer to try to improve the performance of the model. The result can be seen in Figure 9e(g).

Another possible cause for what happens in Figure is the number of observations of the data set. Since the collection of speeches for Bolsonaro's government was interrupted in the article that primarily studied the set, there are only 140 observations to analyze. This amount can be too small, leading to variability in performance metrics. As the idea of this work is to treat the data originally collected, we do not consider the collection of more data as a solution.

5. Discussion and conclusion

This article contributes both methodologically and substantively to the study of political discourse. The results presented offer a deeper understanding of Bolsonaro's rhetoric, particularly through the lens of populism. By comparing speeches classified as populist and non-populist, significant insights were uncovered.

This work analyzes a selection of potentially populist speeches from the database used by Ricci, Izumi, and Moreira (2021). Simple techniques were used to pre-process the speeches of former President Jair Bolsonaro, present in the original document. Some of the most common machine learning algorithms aimed at classification were employed

to evaluate their respective performances. Analyzing the results obtained, it is possible to conclude that the SVC algorithm achieved the best performance, with an accuracy of 86% and the best values for the three other evaluation metrics, considering the two classes in question.

However, as previously mentioned, this result cannot be considered a success due to the overfitting problem. Furthermore, an analysis of the other metrics highlights the difficulty the model encountered in differentiating between populist and non-populist phrases. Despite achieving a good overall accuracy, this value predominantly reflects the positive results in classifying sentences labeled as populist. However, for non-populist phrases, the model did not perform satisfactorily, with a precision of 0.67, recall of 0.5, and F1-score of 0.57.

This limitation reflects the challenges these techniques face in capturing the deeper meaning of the speeches, making them unsuitable for more abstract and complex studies, such as those in Political Science. Nonetheless, this work serves as an important step in outlining the boundaries of these approaches. On the other hand, the analysis of the speeches using alternative techniques proved very useful in understanding the elements mobilized in the narratives constructed by the former president, confirming the populist elements of defending the “sovereignty of the people” and fostering social polarization. It is, therefore, a highly efficient complementary analytical tool for providing a broader perspective of the texts analyzed.

Additionally, this study reinforces the understanding of Bolsonaro’s early presidency as deeply rooted in populist communication strategies. These findings are critical in understanding how Bolsonaro employed rhetorical strategies to connect with specific voter bases.

This work can still evolve by applying more advanced classification techniques to other, more recent databases of speeches, while also applying the approaches that proved effective in this study. This enables the execution of comparative and panel research, allowing for the monitoring and understanding of the populist phenomenon across countries and political actors. Future research could explore how these strategies evolved throughout his term and how they compare with other populist leaders globally.

The use of machine learning algorithms, particularly Support Vector Classification, demonstrated effectiveness in classifying populist rhetoric with a high accuracy rate. However, it also exposed challenges such as data imbalance, which suggests further research is needed to improve model training and testing with larger datasets.

Referências

- Engesser, S.; Ernst, N.; Esser, F.; Büchel, F. (2017) Populism and social media: how politicians spread a fragmented ideology, *Information, Communication & Society*, 20:8, 1109-1126, DOI: 10.1080/1369118X.2016.1207697.
- Joshi, Y. (2020). A quick guide to text cleaning using the nltk library. *Analytics Vidhya*.

Available at: https://www.analyticsvidhya.com/blog/2020/11/text-cleaning-nltk-library/h2_2.

Kirk, H., Rovira, C. (2017). What the (ideational) study of populism can teach us, and what it can't. *Swiss Political Science Review*, 23(4).

Lust, Ellen; Waldner, David. Theories of Democratic change. Retrived from https://pdf.usaid.gov/pdf_docs/PBAAD635.pdf, 2015.

Moffit, Benjamin. The Global Rise of Populism: Performance, Political Style, and Representation. Stanford University Press, 2016.

Nicolau, J. (2020). *O Brasil dobrou a` direita: Uma radiografia da elei,c~ao de Bolsonaro em 2018*. Rio de Janeiro: Zahar. 141 p.

Pérez-Liñán (2018). A Impeachment or backsliding? Threats to democracy in the twenty-first century, *Revista Brasileira de Ciências Sociais* 33;

Ricci, P., Izumi, M., Moreira, D. (2021). O populismo no Brasil (1985-2019): um velho conceito a partir de uma nova abordagem. *Revista Brasileira de Ciências Sociais*, 36(107), 1-22.: <https://doi.org/10.1590/3610707/2021>. Acesso em: 20 ago. 2024.

Torfi, A., Shirvani, R. A., Keneshloo, Y., Tavvaf, N., Fox, E. (2020). Natural language processing advancements by deep learning: A survey. *ArXiv*, vol. abs/2003.01200.

Appendices

Figure 5: Heatmap of the relationship between each two words in the document considered non-populist.

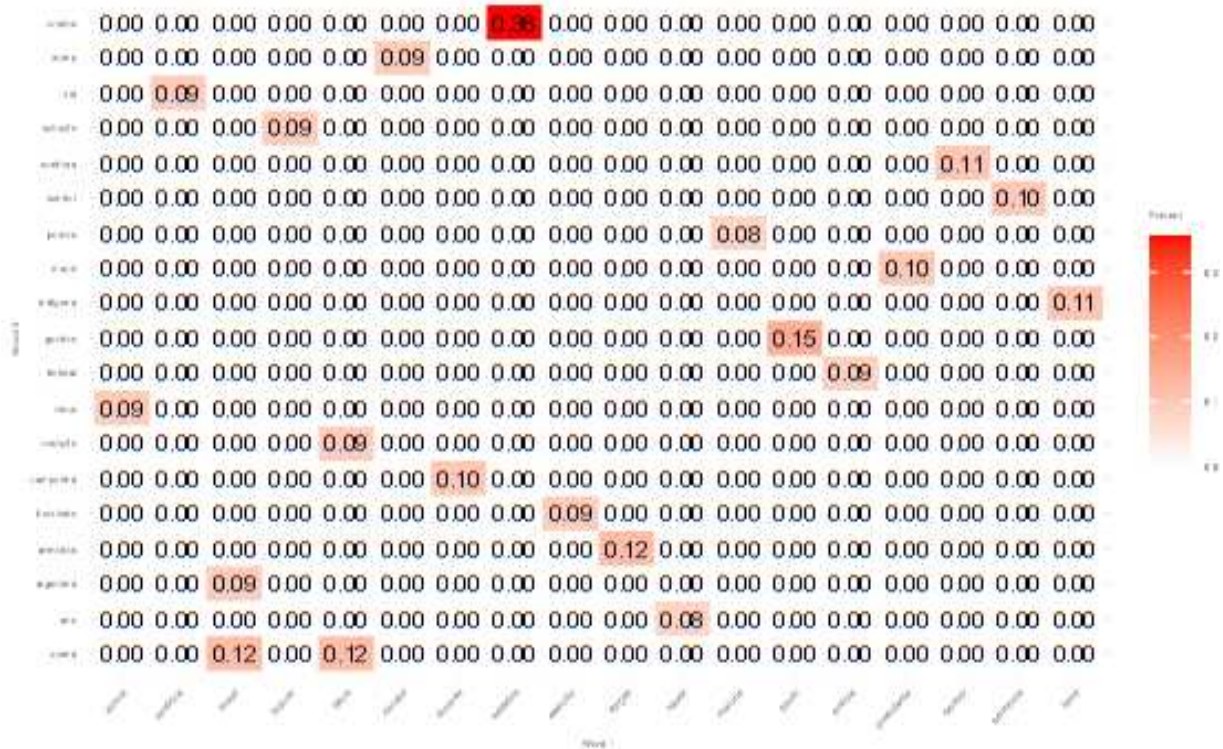


Figure 6: Heatmap of the relationship between each two words in the document considered populist.

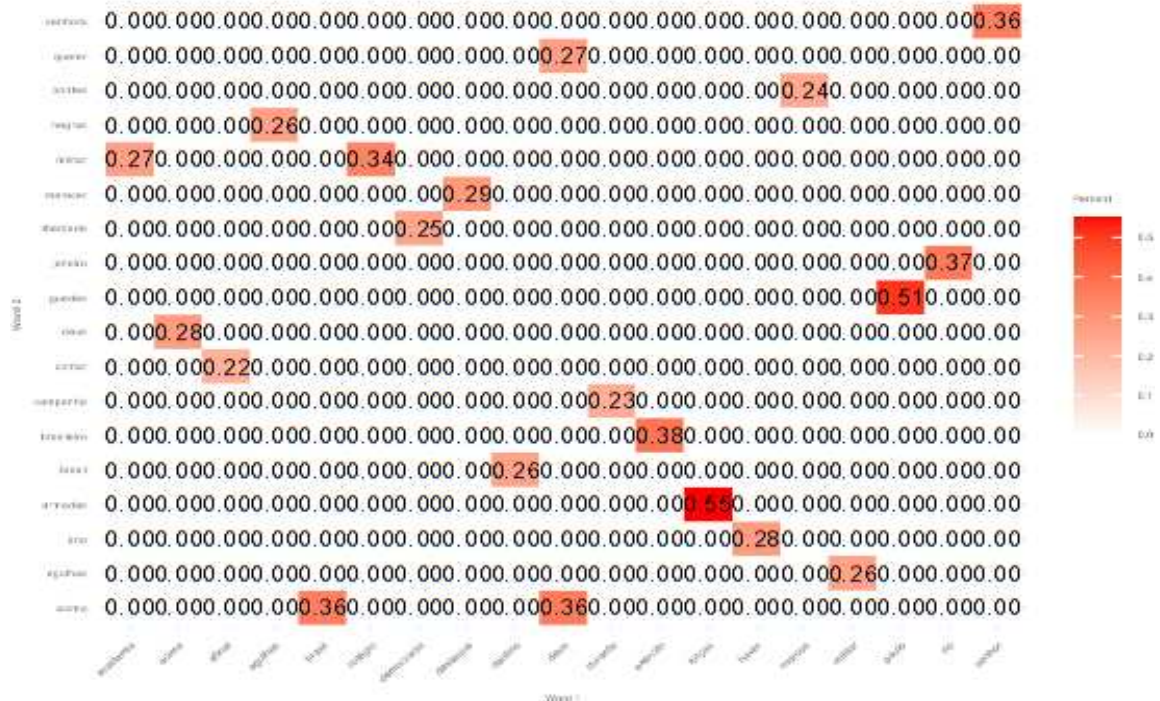


Figure 7: Confusion Matrix

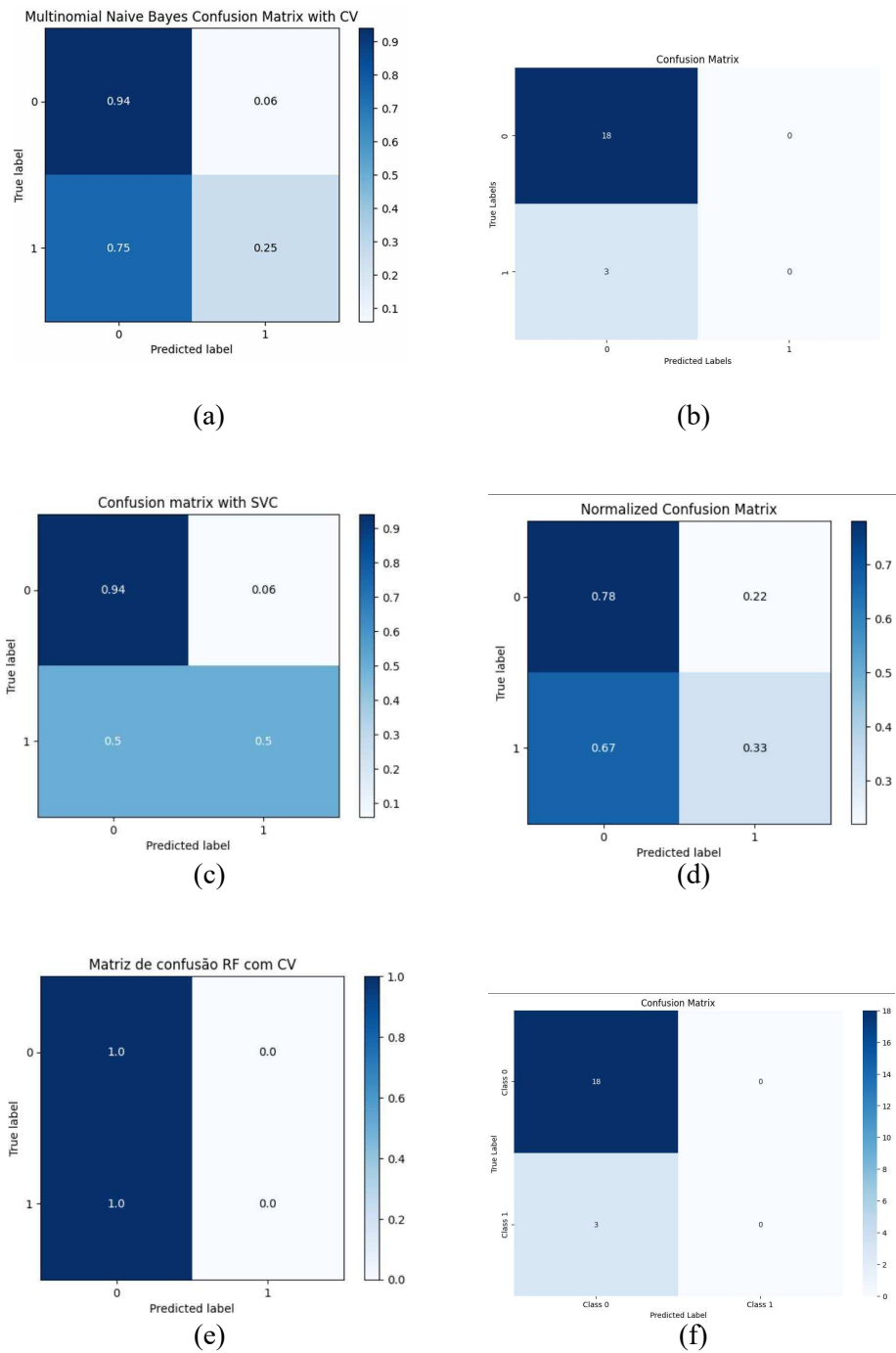
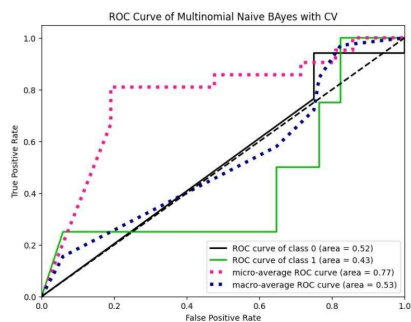
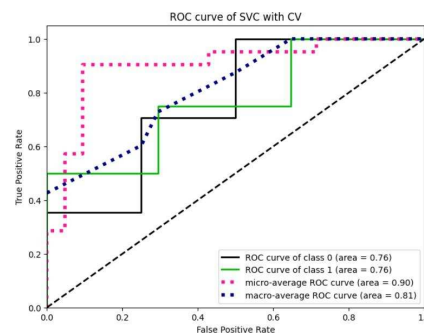


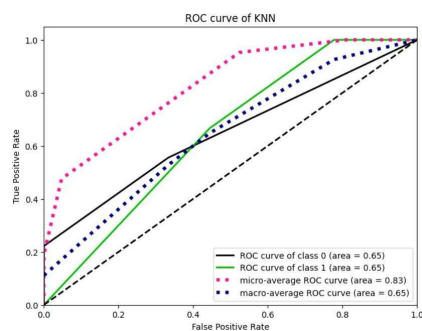
Figure 8: ROC curves



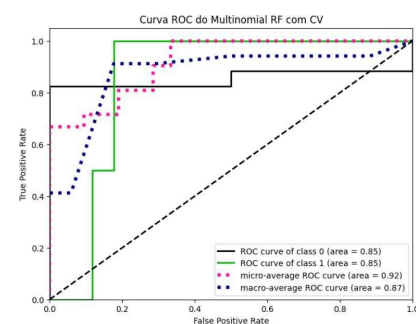
(a)



(b)



(c)



(d)

Figure 9: Learning curves

