# Sentiment Analysis of Portuguese Economic News

## Cátia Tavares ✉ ⌂
Instituto Universitário de Lisboa (ISCTE-IUL), Lisboa, Portugal

## Ricardo Ribeiro ✉ ⌂ 🆔
Instituto Universitário de Lisboa (ISCTE-IUL), Lisboa, Portugal
INESC-ID, Lisbon, Portugal

## Fernando Batista ✉ ⌂ 🆔
Instituto Universitário de Lisboa (ISCTE-IUL), Lisboa, Portugal
INESC-ID, Lisbon, Portugal

---- **Abstract** ----

This paper proposes a rule-based method for automatic polarity detection over economic news texts, which proved suitable for detecting the sentiment in Portuguese economic news. The data used in our experiments consists of 400 manually annotated sentences extracted from economic news, used for evaluation, and about 90 thousand Portuguese economic news, extracted from two well-known Portuguese newspapers, covering the period from 2010 to 2020, that have been used for training our systems. In order to perform sentiment analysis of economic news, we have also tested the adaptation of existing pre-trained modules, and also performed experiments with a set of Machine Learning approaches, and self-training. Experimental results show that our rule-based approach, that uses manually written rules related to the economic context, achieves the best results for automatically detecting the polarity of economic news, largely surpassing the other approaches.

## 1 Introduction

Economic data and economic indicators are an important resource to reveal the true picture of an economy's condition. They allow us to understand the state of the economy and to determine our investment and consumption decisions. Also, forecasters and policy makers need information about how economy stands to take appropriate responses to their decisions. To give response to that and given the fact that economic indicators have usually a lag in their publication, having mostly a monthly and quarterly frequency, we should take advantage of the increasingly digital world that we have and transform the exponential amount of information available in an opportunity. Taking into consideration that news are the main form of transmission of information about the present, they can generate changes in the expectations of their readers. If the news are positive, the expectations of economic agents will also be positive and, consequently, the sentiment about the future of the economy as well. Otherwise, mistrust will be generated about the situation of the economy, which may have repercussions on economic agents investment and consumption actions [22].

Currently, a large amount of information is shared in news sites, blogs, and social networks. If processed timely and adequately, it can help to obtain key insights about the economic situation in almost real time. Sentiment analysis, also known as opinion mining, is the task

of finding the opinion of an author of a text concerning its content [9]. In general, it consists of identifying the polarity, positive or negative, of a text span (e.g., a document or a sentence) that contains "explicit opinions, beliefs, and views about specific entities" (a subjective text span).

Considering that the available linguistic resources for the Portuguese language are scarce, sentiment analysis of texts in Portuguese is still an active field of research, especially considering documents in specific domains, such as the economic domain [23]. In this work, we explore three approaches to perform sentiment analysis of Portuguese economic news. Our baseline approach consists of translating Portuguese economic news data into English, and then apply well-known and widely used sentiment resources for English, such as VADER [15] and TextBlob[1]. In the second approach, we manually created a set of rules based on the economic context, and used a rule-based approach. Finally, we trained different machine learning models, in order to try to improve our results even further. The performance of each one of the previous approaches was evaluated using a manually labeled dataset, containing 400 economic sentences, created in the scope of this work.

The rest of this document is organized as follows: Section 2 presents an overview on the related literature; Section 2.1 overviews the different strategies commonly used to perform sentiment analysis, while Section 2.2 focus on sentiment analysis in the economic context; Section 3 presents a description of our news data in Section 3.1 and the description of our Golden Data manually annotated in Section 3.2; Section 4 presents the proposed pipeline of our work; Section 5 describes the experiments performed with each one of our adopted approaches, namely, the translation-based approach (Section 5.1), the rule-based approach (Section 5.2), and the Machine Learning approach (Section 5.3), with Section 5.4 presenting a summary of the results attained; finally, Section 6 presents the major conclusions, and recommendations for future work.

## 2    Related Work

This section describes some relevant literature. We start by addressing sentiment analysis in general, and then we focus on its application in the economic domain.

## 2.1    Sentiment Analysis

As previously mentioned, sentiment analysis focuses on the analysis of users' expressions, classifying them according to the polarity. Data from different types of sources such as blogs, news, and social media, the use of different languages, non-standard words and the use of emojis and other symbols led to approaches with distinct complexity levels.

Sentiment analysis has gained an important role in the analysis and understanding of consumer communication in the media, allowing to provide key information about the public opinion on several subjects [31]. In that sense, in recent years, in addition to the more traditional focus on data from different news services, research in the field of sentiment analysis has been carried out in several domains, focusing on the analysis of data from social networks such as Twitter [12]. The difference between these two types of textual data is that in the latter the opinion is generally clear, objective and is well defined in the text, while the first may cover several domains and may consist of more subjective texts and descriptions of complex and context-based events [3].

---

[1] `https://textblob.readthedocs.io`

There are several techniques for sentiment analysis. There are approaches based on lexicons, which consist of predefined collections/dictionaries of terms and the associated sentiment/emotion; approaches based on Machine Learning (ML); and, even hybrid approaches, in which both the previous approaches are combined. ML techniques can also be divided into two groups, supervised techniques and unsupervised techniques. In the first case, the data must be labeled, which is not the case with unsupervised techniques [18, 7].

Sentiment dictionaries are dictionaries where each word is associated with and opinion/polarity (positive, negative or neutral) and are very useful resources to classify the sentiment polarity. There are many sentiment dictionaries based on the English lexicon, however, for the Portuguese language these resources are scarce. The literature about sentiment analysis focusing on the English language is vast but the linguistic resources available for sentiment analysis in Portuguese and other languages are still limited. Several studies adopt an approach based on the translation of the original data to English and after that an English sentiment analysis tool is applied. However, translation errors and language specific information can have a significant impact on the final result [23].

In the specific case of Portuguese, well-know sentiment analysis tools, like VADER, TextBlob, or Stanza [24], do not work. VADER combines a lexicon and a rule-based approach for sentiment analysis. VADER original experiments were performed only on English data. TextBlob is a Python library that provides several natural language processing modules, including one for sentiment analysis. TextBlob includes two sentiment analysis approaches, a rule-based model and a supervised ML model, based on a Naïve Bayes classifier. As provided, it only deals with the English language. Stanza toolkit also uses a ML model for sentiment analysis, in this case based on a Convolutional Neural Network classifier. Stanza is also a Python library and has models for English, Chinese, and German.

The ML approaches rely on ML algorithms to determine the sentiment as a text classification problem. Given a phrase/instance of unknown class, the model predicts the label/class to which it belongs. Supervised methods require the existence of labels in the training data, and example of supervised classifiers are Support Vector Machine (SVM) and decision trees [18, 29]. The unsupervised methods do not require the existence of labels and have been subject of a lot of research.

Depending on the approach used to perform SA, before classifying the sentiment it is necessary to extract and select the features of the text. Feature extraction consists of performing a transformation to the original features and generate more significant ones, aiming to reduce complexity and provide simpler representations to the data [16, 18]. During the feature extraction process, useful features are identified and extracted, and analysis can also be done to understand which features increase accuracy the most. To help weight features, measures such as TD-IDF (term frequency-inverse document frequency) are used. After extracting the features, the sentiment classification is done.

Ahmed and Ahmed [2] used an approach based on lexicons to classify data from news. Firstly, they used text-preprocessing techniques, such as punctuation removal, "stop-words" removal and stemming. After reducing the derived words, they computed the polarity using TF-IDF measure to identify the most frequent words and to be able to assign them sentiment scores through dictionaries such as SentiWordNet. Finally, the news polarity was determined as positive, negative or neutral, by the sentiment average of the total news words.

Mohamed [19] evaluated several algorithms to perform sentiment analysis based, concluding that SVM outperforms other methods such as Naïve Bayes and decision trees. However, each SA technique will have different performance and results depending on the data in which it is applied.

## 2.2   News, Sentiment Analysis and the Economy

The relationship between economic news, the economy and public perception, and opinion has been a subject of research for a while. The news have influence on the evaluations and opinions of economic agents about the economy. When they are negative, public opinion about future economic conditions is unfavorable and pessimism about the economy is generated [8, 13]. The importance of public opinion is due to the fact that changes in expectations about the economic future can be a source of economic fluctuations [11].

In order to understand the current state of the economy, high-frequency information is needed quickly and in real time [26]. This way, economic agents can use a multitude of high-frequency information in order to guide their actions, including news from the media [28].

There are several studies that focus on understanding the behavior of financial markets and stock values based on economic news [30], some of them are shown in Table 1. Mining the news plays an important role in designing strategies to predict market behavior and, based on events and news items, it is possible to predict market prices [17]. When there is pessimism in the media, patterns of falling stock prices and short-term returns are expected, concluding that the news information is useful for making predictions of market return and risk [27, 6].

In relation to the foreign exchange market, text mining is also a promising way to predict exchange rate movements based on the economic events present in the news, bringing benefits to investors and risk managers [20].

According to Huang et al. [14], traditional economic indicators based on surveys have been replaced by techniques for extracting sentiments from news texts and central bank statements, through the application of machine learning and other computational techniques. News-based sentiment indicators make it possible to predict periods of financial crisis, serving as early indicators of them. Nyman et al. [21], showed that periods of financial crisis can be detected in the news, being preceded by sentiments of anxiety. From the Bank of England news and publications data, they were able to obtain information about episodes of risk and market volatility.

Aguilar et al. [1], when trying to monitor economic activity in Spain by building a sentiment indicator based on the news, found that the developed indicator has advantages over the indicator based on surveys in GDP forecasting and in forecasting the crisis related to COVID-19. Also according to Fraiberger et al. [10], the sentiment indicator based on the sentiment present in the news gives a direct and real-time view of the aggregate sentiment of the current and future state of the economy, correctly portraying fluctuations in GDP, which allows policy makers to react more efficiently to economic conditions.

## 3   Data

This work uses a large dataset of economic news for training our models, and a dataset of labeled news that serves as reference for our evaluation experiments. This section presents the details about the two datasets.

## 3.1   Portuguese Economic News

The data used for training our models consists of economic news produced between January 1$^{st}$, 2010 to December 31$^{st}$, 2020, covering an 11 years time-span (132 months). The corpus was extracted from online news published by *Expresso* and *Público*, two well-known Portuguese newspapers. For each article, we have extracted the date, the headline and the
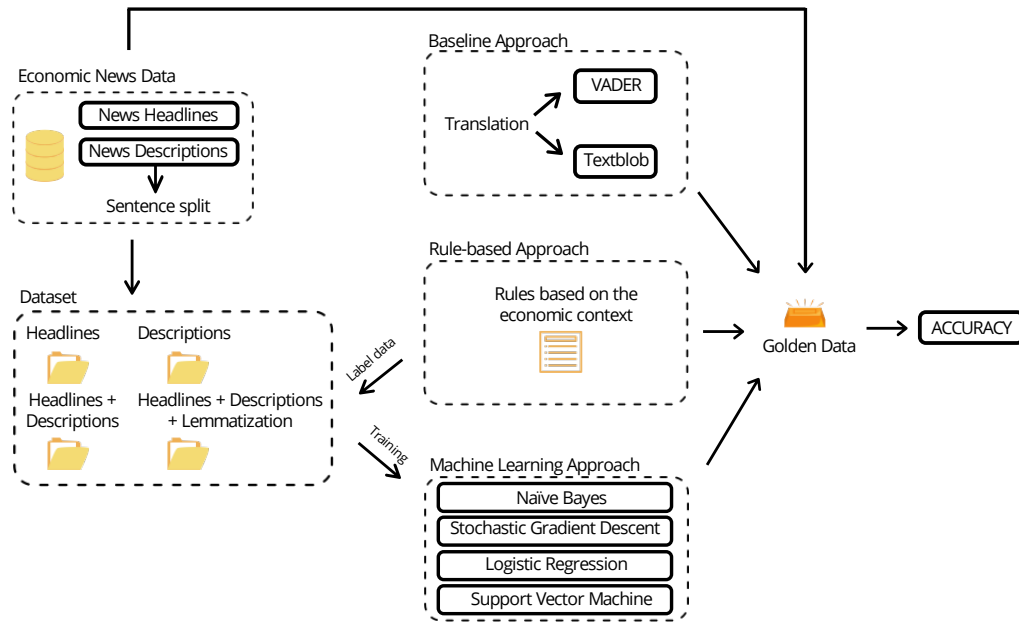
■ **Table 1** Sentiment analysis in economic context.

| Ref. | Goals and Results | Data | Techniques | Metrics |
|------|-------------------|------|------------|---------|
| [14] | Use news data to predict periods of financial crisis<br><br>News sentiment index contains useful information to predict financial crises and market risk | Financial Times News | – Word vector representation<br>– Semantic clustering<br>– Sentiment of each cluster | Precision Recall F-score |
| [21] | Use text analysis to extract statistics about the economy, predicting important events and systemic risk<br><br>Strong correlation with financial market events, such as structural breaks, and with other market measures such as sentiment, confidence, market volatility and systemic risk | – Comments about the bank of England market<br>– Financial market research reports<br>– Economic news | – Word count<br>– Loughran and McDonald sentiment dictionary | Granger causality p-value |
| [5] | Create indicator to predict the state and evolution of the economy in France (GDP)<br><br>Media are a promising tool for economic analysis and have made it possible to forecast French GDP | Le Monde News | – Construction of a sentiment dictionary<br>– Logistic regression | RMSFE |
| [10] | Create sentiment index to predict economic fluctuations<br><br>Index gives a direct and real-time view of the aggregate sentiment of the current and future state of the economy, correctly portraying GDP fluctuations | Economic News | Loughran and McDonald sentiment dictionary (economy) and Young and Soroka dictionary (economy and politics) | Granger causality p-value RMSE |
| [1] | Create sentiment indicator to monitor economic activity in Spain in real time<br><br>Correlation of the indicator developed with the Economic Sentiment Indicator (ESI) of 0.8. Better performance than ESI in forecasting GDP and the economic crisis related to COVID-19. Better GDP forecast when we use the indicator developed compared to the ESI. | Economic News | Count words related to improvements and economic downturns | RMSE |

**Table 2** Analysis of the number of words in the headlines and text descriptions.

|  | Headline | Description |
|---|---|---|
| Average length | 9.08 | 15.31 |
| Median length | 9 | 15 |
| Maximum length | 31 | 235 |



**Figure 1** Automatic classification of economic news pipeline.

corresponding description, when available. This accounts to over 90,000 economic news headlines, 62,326 of them also complemented with a textual description, which can contain one or more sentences. Table 2 shows some statistics about the number of words present in each one of the fields for each article.

## 3.2 Manually Annotated Data

In order to evaluate our approaches, we have collected a sample of 400 sentences from recent economic news, and manually classified them with one of three possible labels, according to its corresponding polarity: Negative, Neutral, and Positive.

## 4 Pipeline

In order to automatically classify the sentiment of Portuguese economic news, we have adopted the following strategy, represented in Figure 1, that consists of a data collection stage, an automatic classification stage, and the evaluation. We have started by collecting the data from online newspapers, as described in Section 3. Each one of the news stories was processed in order to extract the corresponding date, title (headline), and description. Additionally to collecting the data from 2010 to 2020, we also have selected 400 sentences, extracted from most recent news, that were manually annotated with the purpose of evaluating the approaches under study.

In terms of available natural language processing tools, the Portuguese language may be considered a low-resource language, and during the course of this work, we could not find a sentiment analysis tool that could be directly applied to detect the sentiment of a sentence in Portuguese. For that reason, our initial strategy, represented in the top-middle box of Figure 1, consisted of translating the Portuguese sentences into English, and then using one of the existing English tools. However, sentiment analysis is known to be domain dependent, and soon we have realised that the commonly used tools could not be easily applied to the economic news domain. So, in order to overcome this problem, we have manually created a set of rules adapted to the economic domain, and we adopted the rule-based approach represented in the middle box of the figure. Finally, we have used our ruled-based approach to label our large dataset of economic news, and, in order to improve our results even further, we have trained several machine learning models, both in a supervised and semi-supervised way, as represented in the bottom rectangle of Figure 1. All the described approaches are evaluated using the same manually labeled dataset, described in Section 3.

## 5 Experiments and Evaluation

In this section, we present the details about the three different approaches used to perform sentiment analysis. As previously mentioned, first we have tried to use existing tools to perform sentiment analysis, but we soon realised that the resources available for the Portuguese language are scarce. Thus, as a first approach, we translated our data into English and then used VADER and TextBlob to perform the analysis. We concluded that this approach is limited when applied to the economic context. So, we tried a second approach where we observed the most common patterns appearing in economy news stories and created a set of rules to classify each sentence, which proved to perform well in our data. In order to improve our results even further, we experimented a third approach where we trained different machine learning models. In the end of the section, we present a summary with the results obtained with the mentioned approaches.

### 5.1 Baseline/Translation-based Approach

When facing the lack of tools for a given language, one possible immediate solution is to translate the existing data to another language and then use the available tools for that language. In fact, during the course of this work, we did not find any available tools to perform sentiment analysis in Portuguese. As so, we have adopted TextBlob and NLTK VADER [15], two well-known tools for sentiment analysis, with the latter reported to perform well when applied to the finance domain [25]. So, after translating our reference data from Portuguese to English using Googletrans, a python library that implemented Google Translate API, VADER achieved an accuracy of 46.5% and TextBlob achieved an accuracy of 32.0%. These results show that this approach is not suitable to the economic context, which was not an unexpected result since we know that sentiment analysis is a domain-dependent task.

We have then applied these tools to our unlabeled data in order to analyse the results in more detail. From the analysis we have observed that, for example, headlines with negative words like unemployment, crisis, deficit, etc., were classified incorrectly most of the times. In fact, the polarity associated with these words is negative, although we have seen that many news involving these words, such as "*unemployment is decreasing*" and "*crisis is slowing down*", should be positive, and that VADER and TextBlob were not taking that into consideration.

■ **Table 3** Expressions related to "unemployment".

| Word 1 | Word 2 | Word 3 | Sentiment |
|---|---|---|---|
| unemployment | reaches | minimum | 1 |
| unemployment | reaches | maximum | −1 |
| unemployment | decreased | | 1 |
| unemployment | increased | | −1 |
| ... | | | |
| unemployment | | | −1 |

## 5.2   Rule-based Approach

Our rule-based approach is similar to the approach proposed by Aguilar et al. [1] for classifying the polarity of economic news, where the sentiment attributed to each headline is also based on rules. For example, when combined with the word "economy", the word "increase" becomes positive, and the word "decrease" becomes negative.

After identifying the errors and limitations in the classification performed by the previous approach, we have started looking at the more prominent words and combination of words in our unlabeled data. We have manually analysed through frequency analysis, the set of words co-ocurring with words like "unemployment", to understand the most frequent patterns. As a result of the analysis of these associations of words, we constructed a list of expressions/rules related to the economic context, and labeled the sentiment associated to them. We have observed meaningful combinations of two and three words, which derived in rules of one, two and three words, accordingly. For example, for the word "unemployment", we can think of expressions involving words such as the ones presented in Table 3. We did the same for other words related to the economic context such as "consumption", "debt", "economy", "recovery", and we end up with approximately 600 rules with the corresponding associated sentiment (-1 if the expression is negative and 1 if it is positive).

Algorithm 1 details our rule-based classification process. First we try to match all the rules involving 3-words expressions. If more than one rule can be applied, we sum the sentiment associated with all the matching rules, and check if the resulting sum is positive, negative, or neutral. If none of the 3-words rules matches our sentence, we try to search all the rules involving 2-words, and again sum the sentiment of each one that we find. Finally, if none of the 2-words rules matches our sentence, we try to match 1-word rules. At the end of this process, if the sentence did not match any rule, then we assume it is neutral. When applied to our corpus of 90,000 news, 9.5% of the headlines where classified as negative and 7.5% as positive. The descriptions, 5.8% were classified as negative and 5.8% as positive. When applied to our reference data, our rule-based approach achieves an accuracy of 86.3%, a significant improvement over the baseline approaches.

## 5.3   Machine Learning Approach

In order to improve the results even further, we have performed additional experiments using our unlabeled economic news dataset for training our machine learning models.

The dataset described in Section 3.1 was used to create four different collections of texts that were used in our Machine Learning experiments: 1) sentences extracted from the headlines (about 90,000 sentences); 2) sentences extracted from the descriptions (about 85,000 sentences); 3) a combination of the previous two collections (about 175,000 sentences); and 4) a combination of the previous collection with its variant where lemmatization was

■ **Algorithm 1** Classification of each sentence based in the rules created.

```
input: sentence, rules

sentiment = 0
for rule in [rules with 3 words]:
    if rule.applies_to(sentence):
        sentiment += rule.sentiment
if sentiment = 0 then
    for rule in [rules with 2 words]:
        if rule.applies_to(sentence):
            sentiment += rule.sentiment
if sentiment = 0 then
    for rule in [rules with 1 word]:
        if rule.applies_to(sentence):
            sentiment += rule.sentiment

if sentiment < 0 then
    return -1
else if sentiment > 0 then
    return 1
else
    return 0
```

applied to the words in the texts, in order to capture a broader set of economic terms (about 350,000 sentences). Lemmatization is a preprocessing step often used in text mining and natural language processing, that consists of converting each word in its basic/root form, analyzing its morphology in order to remove the inflected affixes, leaving only the lemma [4]. For this task we used Spacy and, for example, this process will convert the words "increasing", "increased", and "increases" into the word "increase".

Our rule-based classifier was used to classify all the sentences in each one of the collections. Then, we have converted all the labeled sentences into their corresponding document representation, using unigrams, bigrams, and trigrams.

We have applied the following classical supervised machine learning methods, used extensively for classification and regression tasks: Naïve Bayes (NB), Stochastic Gradient Descent (SGD), Logistic Regression (LR), and Support Vector Machines (SVM). Each one of the methods was applied to each one of the four previously described text collections, using their default parameters. Concerning the feature weights, we have used simple counts for Naïve Bayes, and Term Frequency – Inverse Document Frequency (TF-IDF) weights for all the other methods. The corresponding evaluation results for our reference data are presented in Table 4.

The results attained show that, in general, the text contained in the title is better for training than the text of the descriptions, but the best results are achieved when combining both fields. With NB and SVM we could see that the use of lemmatization contributed to a better result, we could not see the same when using SGD and LR.

Using only the title texts for training leads to a better accuracy with LR (74.0%), and using only the sentences from descriptions perform better with SGC (73.0%). The collection containing the titles and the descriptions led to better accuracy scores for SGD (76.3%), but the best accuracy (77.0%) was achieved by training the SVM with features produced with lemmatization.

**Table 4** Model evaluation in our reference Data.

| Model | Accuracy | Precision | Recall | F1 |
|---|---|---|---|---|
| VADER (baseline) | 0.465 | 0.477 | 0.465 | 0.470 |
| TextBlob (baseline) | 0.320 | 0.352 | 0.320 | 0.301 |
| Rule-based approach | **0.863** | **0.863** | **0.863** | **0.863** |
| Naïve Bayes | | | | |
|    Titles | 0.615 | 0.638 | 0.615 | 0.607 |
|    Descriptions | 0.593 | 0.645 | 0.593 | 0.591 |
|    Titles + Descriptions | 0.633 | 0.648 | 0.634 | 0.627 |
|    Titles + Descriptions + Lemmatization | 0.663 | 0.675 | 0.663 | 0.661 |
| Stochastic Gradient Descent | | | | |
|    Titles | 0.740 | 0.739 | 0.740 | 0.739 |
|    Descriptions | 0.730 | 0.743 | 0.730 | 0.733 |
|    Titles + Descriptions | 0.763 | 0.763 | 0.763 | 0.762 |
|    Titles + Descriptions + Lemmatization | 0.740 | 0.738 | 0.740 | 0.739 |
| Logistic Regression | | | | |
|    Titles | 0.738 | 0.737 | 0.738 | 0.737 |
|    Descriptions | 0.693 | 0.718 | 0.693 | 0.698 |
|    Titles + Descriptions | 0.758 | 0.764 | 0.758 | 0.759 |
|    Titles + Descriptions + Lemmatization | 0.758 | 0.764 | 0.758 | 0.759 |
| Support Vector Machine | | | | |
|    Titles | 0.738 | 0.736 | 0.738 | 0.736 |
|    Descriptions | 0.678 | 0.697 | 0.678 | 0.683 |
|    Titles + Descriptions | 0.755 | 0.761 | 0.755 | 0.757 |
|    Titles + Descriptions + Lemmatization | 0.770 | 0.773 | 0.770 | 0.771 |

We also have performed additional self-training classification experiments, considering all the texts labeled as positive or negative as the initial labels, and performing label propagation to all the remainder data. Nonetheless, the results achieved did not surpass our previous reported results.

## 5.4    Summary

Table 4 summarizes the results achieved with each one of the approaches when evaluated using our reference data. The baseline approaches, using the NLTK VADER and TextBlob, performed poorly in the economic context, where accuracies of 46.5% and 32.0% were obtained, respectively. The best result was obtained with the rule-based approach with an accuracy of 86.3%. The machine learning approaches were not able to surpass our rule-based system: the best result of the machine learning approaches was achieved using SVM, with an accuracy of 77.0%.

## 6    Conclusions and Future Work

The lack of tools for sentiment analysis for Portuguese and the difficulty to obtain a labeled dataset to train a sentiment analysis system for economic context led us to explore a set of approaches in order to solved this practical problem. First, we have tried a baseline approach where we translated our texts into English and used well-known sentiment analysis tools,

such as NLTK VADER and TextBlob. Given the poor results achieved in the economic context, we tried a rule-based approach for which we have created manual rules, based on the economic domain, and used those rules to classify the polarity of each economic text. Finally, we have created a set of machine learning models, based on the large amount of economic texts that we had available, aiming at improving our results even further.

In order to compare and evaluate the performance of the proposed approaches, we have also created a reference dataset, containing 400 economic sentences, manually classified. The performed experiments have shown that the baseline approach achieves poor results, when applied to the economic domain. The rule-based approach achieved an impressive performance of 86.3% accuracy, a significant increase of performance over the baseline approaches. The machine learning models that we have explored were not able to generalise and surpass the rule-based approach.

Our rule-based approach still lacks proper treatment of the negation, and adversative conjunctions. In the near future, in addition to the rules created, we plan to improve the classifier by treating differently words after a word such as "not" or "don't", and consider ways of dealing with the classification of sentences with adversative conjunctions. Concerning the negation, we should classify each sentence with the opposite sentiment of the rule it matches, for example, "*unemployment did not increase*" should have a positive polarity, whereas "*unemployment increased*" has a negative one. Adversative conjunctions introduce additional challenges. They express opposition or contrast between two statements and it is difficult even for a human, to tell the sentiment that it expresses. For example, in "*unemployment decreased but GDP increased*", the statement before the conjunction "*but*" has a positive sentiment, but the statement after it has a negative sentiment. Our rule-based approach would assign a neutral sentiment to this example, since it matches both negative and positive rules.

### References

1    Pablo Aguilar, Corinna Ghirelli, Matías Pacce, and Alberto Urtasun. Can News Help Measure Economic Sentiment? An Application in COVID-19 Times. *SSRN Electronic Journal*, 2020. `doi:10.2139/ssrn.3673825`.

2    Jeelani Ahmed and Muqeem Ahmed. A framework for sentiment analysis of online news articles. *International Journal on Emerging Technologies*, 11(3):267–274, 2020.

3    Alexandra Balahur, Ralf Steinberger, Mijail Kabadjov, Vanni Zavarella, Erik Van Der Goot, Matina Halkia, Bruno Pouliquen, and Jenya Belyaeva. Sentiment analysis in the news. In *Proceedings of the 7th International Conference on Language Resources and Evaluation, LREC 2010*, pages 2216–2220. ELRA, 2010. `arXiv:1309.6202`.

4    Ivan Boban, Alen Doko, and Sven Gotovac. Sentence retrieval using Stemming and Lemmatization with different length of the queries. *Advances in Science, Technology and Engineering Systems*, 5(3):349–354, 2020. `doi:10.25046/aj050345`.

5    Clément Bortoli, Stéphanie Combes, and Thomas Renault. Nowcasting GDP Growth by Reading the Newspapers. *Economie et Statistique / Economics and Statistics*, 505-506:17–33, 2018. URL: `https://EconPapers.repec.org/RePEc:nse:ecosta:ecostat_2018_505-506_2`.

6    Charles W. Calomiris and Harry Mamaysky. How news and its context drive risk and returns around the world. *Journal of Financial Economics*, 133(2):299–336, 2019. `doi:10.1016/j.jfineco.2018.11.009`.

7    Cagatay Catal and Mehmet Nangir. A sentiment classification model based on multiple classifiers. *Applied Soft Computing Journal*, 50:135–141, 2017. `doi:10.1016/j.asoc.2016.11.022`.

**8**     Alyt Damstra and Mark Boukes. The Economy, the News, and the Public: A Longitudinal Study of the Impact of Economic News on Economic Evaluations and Expectations. *Communication Research*, page 009365021775097, 2018. `doi:10.1177/0093650217750971`.

**9**     Ronen Feldman. Techniques and applications for sentiment analysis. *Commun. ACM*, 56(4):82–89, 2013. `doi:10.1145/2436256.2436274`.

**10**    Samuel P. Fraiberger. News Sentiment and Cross-Country Fluctuations. *SSRN Electronic Journal*, pages 1–18, 2016. `doi:10.2139/ssrn.2730429`.

**11**    Ippei Fujiwara, Yasuo Hirose, and Mototsugu Shintani. Can News Be a Major Source of Aggregate Fluctuations? A Bayesian DSGE Approach. *Journal of Money, Credit and Banking*, 43(1):1–29, 2011. `doi:10.1111/j.1538-4616.2010.00363.x`.

**12**    Elena Georgiadou, Spyros Angelopoulos, and Helen Drake. Big data analytics and international negotiations: Sentiment analysis of Brexit negotiating outcomes. *International Journal of Information Management*, 51(November):102048, 2020. `doi:10.1016/j.ijinfomgt.2019.102048`.

**13**    Joe Bob Hester and Rhonda Gibson. The economy and second-level agenda setting: A time-series analysis of econom... *Journalism & Mass Communication Quarterly*, 80(1), 2003.

**14**    Chengyu Huang, Sean Simpson, Daria Ulybina, and Agustin Roitman. News-based Sentiment Indicators. *IMF Working Papers*, 19(273), 2019. `doi:10.5089/9781513518374.001`.

**15**    Clayton J. Hutto and Eric Gilbert. VADER: A parsimonious rule-based model for sentiment analysis of social media text. In Eytan Adar, Paul Resnick, Munmun De Choudhury, Bernie Hogan, and Alice H. Oh, editors, *Proceedings of the Eighth International Conference on Weblogs and Social Media, ICWSM 2014, Ann Arbor, Michigan, USA, June 1-4, 2014*. The AAAI Press, 2014. URL: `http://www.aaai.org/ocs/index.php/ICWSM/ICWSM14/paper/view/8109`.

**16**    Samina Khalid, Tehmina Khalil, and Shamila Nasreen. A survey of feature selection and feature extraction techniques in machine learning. *Proceedings of 2014 Science and Information Conference, SAI 2014*, pages 372–378, 2014. `doi:10.1109/SAI.2014.6918213`.

**17**    Anuj Mahajan, Lipika Dey, and Sk Mirajul Haque. Mining financial news for major events and their impacts on the market. In *Proceedings - 2008 IEEE/WIC/ACM International Conference on Web Intelligence, WI 2008*, pages 423–426, 2008. `doi:10.1109/WIIAT.2008.309`.

**18**    Walaa Medhat, Ahmed Hassan, and Hoda Korashy. Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4):1093–1113, 2014. `doi:10.1016/j.asej.2014.04.011`.

**19**    Ayman Mohamed. An Evaluation of Sentiment Analysis and Classification Algorithms for Arabic Textual Data. *International Journal of Computer Applications*, 158(3):29–36, 2017. `doi:10.5120/ijca2017912770`.

**20**    Hamed Naderi Semiromi, Stefan Lessmann, and Wiebke Peters. News will tell: Forecasting foreign exchange rates based on news story events in the economy calendar. *North American Journal of Economics and Finance*, 52(December 2018):101181, 2020. `doi:10.1016/j.najef.2020.101181`.

**21**    Rickard Nyman, Sujit Kapadia, and David Tuckett. News and narratives in financial systems: exploiting big data for systemic risk assessment. *Journal of Economic Dynamics and Control*, 2021. `doi:10.1016/j.jedc.2021.104119`.

**22**    Nataliia Ostapenko. *Macroeconomic Expectations: News Sentiment Analysis*. Working Paper Series. Eesti Pank, 2020. `doi:10.23656/25045520/052020/0178`.

**23**    Denilson Alves Pereira. A survey of sentiment analysis in the Portuguese language. *Artificial Intelligence Review*, 54(2):1087–1115, 2020. `doi:10.1007/s10462-020-09870-1`.

**24**    Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020. URL: `https://nlp.stanford.edu/pubs/qi2020stanza.pdf`.

**25**    Sahar Sohangir, Nicholas Petty, and DIngding Wang. Financial Sentiment Lexicon Analysis. In *Proceedings - 12th IEEE International Conference on Semantic Computing, ICSC 2018*, volume 2018-January, pages 286–289. Institute of Electrical and Electronics Engineers Inc., April 2018. `doi:10.1109/ICSC.2018.00052`.

**26**  Michael Stanger. A Monthly Indicator of Economic Growth for Low Income Countries. *IMF Working Papers*, 20(13), 2020. `doi:10.5089/9781513525853.001`.

**27**  Paul C. Tetlock. Giving content to investor sentiment: The role of media in the stock market. *Journal of Finance*, 62(3):1139–1168, 2007. `doi:10.1111/j.1540-6261.2007.01232.x`.

**28**  Leif Anders Thorsrud. Nowcasting Using News Topics. Big Data versus Big Bank. In *Norges Bank Research*, Working papers from Norges Bank. Norges Bank, 2017. `doi:10.2139/ssrn.2901450`.

**29**  Yenny Villuendas-Rey, Carmen F. Rey-Benguría, Ángel Ferreira-Santiago, Oscar Camacho-Nieto, and Cornelio Yáñez-Márquez. The Naïve Associative Classifier (NAC): A novel, simple, transparent, and accurate classification model evaluated on financial data. *Neurocomputing*, 265:105–115, 2017. `doi:10.1016/j.neucom.2017.03.085`.

**30**  Ritu Yadav, A. Vinay Kumar, and Ashwani Kumar. News-based supervised sentiment analysis for prediction of futures buying behaviour. *IIMB Management Review*, 31(2):157–166, 2019. `doi:10.1016/j.iimb.2019.03.006`.

**31**  Shanshan Yi and Xiaofang Liu. Machine learning based customer sentiment analysis for recommending shoppers, shops based on customers' review. *Complex & Intelligent Systems*, 6(3):621–634, 2020. `doi:10.1007/s40747-020-00155-2`.