

# Topic categorization in Portuguese news articles

André F. Santos

*CRACS & INESC-Porto LA*  
*Faculty of Sciences, University of Porto*  
Porto, Portugal  
afs@inesctec.pt

**Abstract.** Categorizing news articles according to their contents allows to decrease the information entropy in a world where the rate of publication of digital text documents is increasing fast. In this article we describe ongoing work which aims to evaluate the feasibility of implementing a classifier which is lightweight enough to be used in real time on the client side of a web application. More specifically, we gathered a corpus of Portuguese news and used it to train and evaluate several classification algorithms. We analyze the results obtained in terms of the classifiers error rate, training time and memory footprint.

**Keywords:** topic categorization, machine learning, text mining

## 1 Introduction

Online news articles first appeared as reprints from traditional newspapers; nowadays, however, they represent now the primary source of news for some segments of the population, both in developed and developing countries (whether consumed directly in the newspaper website, or indirectly e.g. through a social media application or a feed catcher) [2, 6, 8].

Unofficially known as *the fourth branch of government*, the press plays a vital role within our society, keeping us informed about the current state of affairs (at a local and global scale) and acting as a watchdog for the other three branches (legislative, executive and judicial). The (lack of) freedom of press and access to the news in a given country is even often considered an indicator of a lack of democracy [7, 10].

As such, improving the ways citizens can access the information (view it, query it and search for it) contained in news articles has the potential to contribute for a more informed and, ultimately, better, society [3].

On the other hand, the last decades have witnessed a fast increase on the rate of publication of digital text documents. Traditional document types, such as news articles, scientific papers or books are now published online along with new formats, such as blog posts or tweets, each having thousands or millions of new documents published each day [1, 9].

Publication is not the only step which has moved to the digital world; in fact, most often nowadays the whole document lifecycle happens digitally, with virtual tools available for researching, writing, styling, publishing and sharing [17].

Having the entire workflow happening within the digital world presents some opportunities when compared to the traditional process [12]. In particular, due to the current processing power commonly available, tasks related to the manipulation of the information contained within these documents (searching, compiling, annotating, sharing, ...) can now be performed automatically and targeting a large amount of articles.

In addition to the document content (for example, in a news article, the *title*, *lead* and *body*), its metadata is also important: author(s), date of publication, source, topic, mentioned entities and their relations, etc [18,19]. Some of this metadata might be filled in and stored along with the document (e.g. *author* and *date of publication*); other is usually extracted from the document content (e.g. mentioned entities) [16].

An example of a feature which improves information access is the categorization of news articles by the topic (or topics) of its content [11]. The presence of such a categorization may influence the way the information is stored, organized, displayed and queried [15].

The image shows a web form interface for classifying news articles. At the top is a search bar. Below it is a 'Title' field. The 'Body' section contains several lines of placeholder text. The 'Categories' section features a dropdown menu with 'International' selected. Below the dropdown, there are two buttons labeled 'Politics' and 'Economics' under the heading 'Suggestions'.

**Fig. 1.** Category classification and suggestion on the client side

The simplest way of achieving this categorization is to have the author of the article manually introducing it (e.g. the journalist typing it on the news article authoring framework); however, this solution presents some challenges:

- It increases the amount of work the author has to do.
- The author might not be sure which categories are available.
- The author might not be sure which category is the best (e.g. *Economics* vs *Finance*).
- It does not scale – e.g. if the goal is to categorize an existing (large) corpus.

Thus, an automated way of categorizing news articles could solve some of these problems and decrease the burden of this task. Additionally, a lightweight version of such a classifier could be implemented on the client side code of a web application, for example, allowing the categorization to happen in real time (i.e. as the author types in the article text). Figure 1 presents a suggestion of how this feature could look like if implemented on a web application.

The challenges of document classification have been well studied within the machine learning research field of study [4, 13, 14]. Given a corpus of already classified documents, several algorithms might be applied to train a classifier capable of determining the category of additional articles.

In this article, we describe the preliminary results obtained in developing a classifier to categorize news articles using a previously manually categorized corpus. Additionally, we evaluate the possibility of implementing such a classifier as lightweight as possible to allow it to run on the client side of a web application.

## 2 Methods

In order to train and evaluate classification algorithms, we first needed to choose and obtain a suitable dataset. Preferably, this dataset should contain documents which were previously categorized, allowing us to skip the time and effort-consuming task of manually categorizing the articles ourselves. Once this dataset was chosen and obtained, we would then clean and prepare it to be used to train the classifiers.

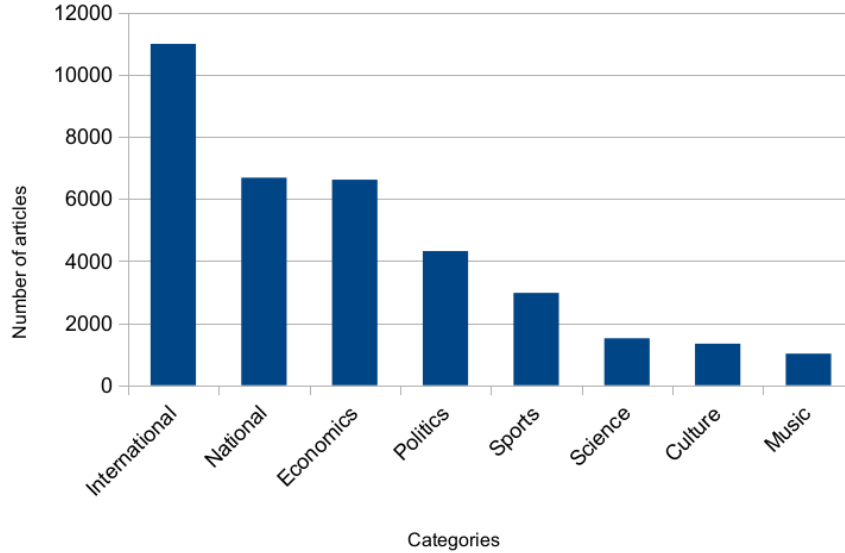
### 2.1 The dataset

We gathered a dataset of news articles published in *Observador*<sup>1</sup>, one of the main Portuguese newspapers, which stands out from the others for being fairly young (it was created in May 2014) and for existing exclusively online. The initial dataset comprised 42.475 entries, the most recent ones dated from November 2016, from which we used only a subset, for reasons later described.

We gathered all the categories used by Observador, and ordered them from the most common to the least common. We selected the ones which had more than 1000 articles in our dataset, and reduced our original dataset to include

<sup>1</sup> <http://observador.pt>

articles from these categories only. Figure 2 presents an overview of the selected categories and the number of articles available for each one.



**Fig. 2.** Total number of articles retrieved for each category

We then randomly selected, from each category, 700 articles to be used to train the classifiers, and 200 to be used to evaluate their performance.

For each article, we had available its contents (title, lead, body) and several metadata fields (publication date, category, tags, etc). A truncated JSON representation of an article can be found in Listing 1.

## 2.2 Preprocessing the articles

Originally, the dataset was obtained as a large MongoDB collection (more than 2.5 million entries), containing articles from several Portuguese and international newspapers. The process needed to transform this collection into data our classifiers could process required querying the database, exporting the news articles and splitting them into a train and an evaluation datasets.

The database query selected articles from *Observador* where the body had a length greater than 100 characters (to discard some malformed articles which had an empty body or a body composed of only a few words), and the categories included at least one of the most common categories.

For each article returned by the query, the pretitle, title, subtitle and lead fields, if present, were simply copied to a plain text file, separated by blank lines. The body field, however, was stored in the database in HTML format. As such,

```
{
  Type: "sapo.obj.creativework.article",
  Source: {
    Name : "Observador"
  },
  Pretitle: "Benfica",
  Title: "Ruben Amorim com rotura total do ligamento cruzado",
  Author: {
    Name: "Observador"
  },
  Tags: [
    "benfica",
    "desporto",
    "futebol",
    "ruben amorim"
  ],
  PublishDate: ISODate("2014-08-25T18:33:00Z"),
  Lead: "Depois de Fejsa, mais uma baixa. O internacional português [...]",
  Body: "<p>O pior cenário confirmou-se. O Benfica informou esta segunda-feira [...]",
  URL: "http://observador.pt/2014/08/25/ruben-amorim-com-rotura-total-ligamento-cruzado/",
  CategoryPaths: ["Desporto"],
  Domain: "observador.pt",
  Language: "pt_PT",
  ...
}
```

**List. 1.** Example of JSON representation of an article

the HTML tags had to be stripped, and then it was also added to the plain text file.

The files were stored in folders, separated by category. For each category, 700 articles were allocated to the train set, and 200 to the evaluation set. These preprocessing tasks were accomplished using Bash and Node.js scripts.

Both datasets were loaded into R using the **tm**<sup>2</sup> package. Each was then passed to a function responsible for preprocessing the text of the articles:

- The text was converted to lowercase characters.
- The Portuguese stopwords were removed.
- Diacritics were converted to their normalized form (e.g. `à` → `a`).
- Punctuation signs were removed.
- Numbers were removed.
- Word were stemmed (e.g. `conseguiram` → `consequi`).
- White space was removed and the text was tokenized.

A document-term matrix was calculated using the **DocumentTermMatrix** method from the **tm** package. For each of the 5600 documents present in the train set the matrix included all the terms which met the following requirements:

- The term length was between 3 and 30 characters (to discard things like URLs or badly tokenized sentences).
- The term appeared in the document at least twice.
- The term appeared at least in 10 documents.

The algorithm used to weigh the terms in the matrix was *tf-idf* [5]. The obtained matrix presented a sparsity of 99% and contained 3234 distinct terms.

The final step was to remove the sparse terms from the matrix. This allows to discard terms which might be too specific of the train set and which might negatively influence the performance of the classifiers by overfitting them to the train set. Additionally, a smaller list of terms reduces the execution time of both training and applying the classifier, and the memory footprint of the classifier. However, while discarding sparse terms we might end up removing relevant terms and increase our classifiers error rate.

We used the **removeSparseTerms** function from the **tm** package, and produced three distinct lists of terms:

- $LT_{90}$  contained terms with 90% or less sparsity.
- $LT_{95}$  contained terms whose sparsity was under 95%.
- $LT_{99}$  contained the terms with a sparsity level below 99%.

### 2.3 Classification

To create the classifiers, we used 5 well known algorithms: decision tree (DT), k-nearest neighbors (KNN), naive Bayes (NB), neural network (NN) and support vector machine – with radial kernel ( $SVM_{RK}$ ) and linear kernel ( $SVM_{LK}$ ).

<sup>2</sup> <https://cran.r-project.org/web/packages/tm/>

We trained each of the classifiers three different times, one for each list of terms ( $LT_{99}$ ,  $LT_{95}$  and  $LT_{90}$ ). Then we evaluated each trained classifier using the test dataset.

The test dataset contained 1600 documents (200 belonging to each category) and was preprocessed in a way similar to the train set (described in the previous section, 2.2). In each iteration, however, the final list of terms in each test document was restricted to terms present in the corresponding list of terms ( $LT_{99}$ ,  $LT_{95}$  and  $LT_{90}$ ).

### 3 Results

A number of measurements and metrics were calculated regarding the lists of terms, the classifiers training process and their results in the evaluation process.

Table 1 presents the size of each list of terms after removing the terms whose sparsity was above the corresponding threshold.

**Table 1.** Lists size

	Sparsity threshold (%)	Number of terms left
$LT_{90}$	90	45
$LT_{95}$	95	139
$LT_{99}$	99	912

Table 2 presents the execution time for training the algorithm with the lowest error rate for each list of terms.

**Table 2.** Execution times for training

	Algorithm	Training time
$LT_{90}$	DT	12s
$LT_{95}$	KNN	3m
$LT_{99}$	KNN	57m

Table 3 presents the error rates obtained for each list of terms using each of the classifiers, with the value of the best classifier for each list highlighted in bold.

**Table 3.** Error rates

	DT	KNN	NB	NN	SVM <sub>RK</sub>	SVM <sub>LK</sub>
$LT_{90}$	<b>0.60</b>	0.71	0.85	0.61	0.88	0.88
$LT_{95}$	0.70	<b>0.56</b>	0.83	0.58	0.87	0.87
$LT_{99}$	0.70	<b>0.35</b>	0.87	0.76	0.87	0.87

Table 4 presents the confusion matrix generated using the k-nearest neighbors classifier with the  $LT_{99}$  list of therms, with the number of correct classifications for each category highlighted in bold.

**Table 4.** Categories confusion matrix ( $LT_{99}$  with KNN)

	science	culture	sports	economics	international	music	national	politics
science	<b>127</b>	19	1	6	12	24	11	0
culture	5	<b>102</b>	2	2	3	79	6	1
sports	1	3	<b>168</b>	7	3	14	4	0
economics	5	3	0	<b>146</b>	5	17	14	10
international	12	12	7	17	<b>90</b>	35	16	11
music	0	2	0	0	2	<b>184</b>	1	0
national	9	9	9	19	15	31	<b>88</b>	21
politics	1	3	0	30	8	17	14	<b>127</b>

## 4 Discussion

The analysis of the results obtained should take into account other metrics besides the error rates obtained for each classifier.

The size of the list of terms, for example, gives us an idea of the memory footprint of a classifier, a parameter which is of the utmost importance if the goal is to implement a classifier as lightweight as possible. The training time is also relevant, as shorter training times give the possibility of retraining the classifiers more often, allowing them to be updated as the corpus of articles grows in size.

Looking at the actual results obtained and represented in Tables 1 and 2 we can see that  $LT_{90}$  has simultaneously the smallest list of terms (45) and the shortest training time (12 seconds using the DT algorithm). However,  $LT_{90}$  also presents the worst error rates, even when looking at the algorithm which achieved its best results (0.60 using a DT classifier).

On the opposite side,  $LT_{99}$  presented the lowest error rates (0.35 using a k-nearest neighbors classifier), but it took almost one hour to train and used a list comprising 912 terms.



The results obtained confirm that there is a tradeoff between the size of the list of terms and the training time, on the one hand, and the classifier error rate, on the other. However, a list of 912 terms seems to be an acceptable memory footprint for a client side classifier; additionally, the higher training time would not present much problems in this scenario, as the classifier would be trained beforehand and thus not be visible on the client side.

Given the reasonable values for the training time and the size of the list of terms, and looking to the error rate obtained in each of tests performed, we can conclude that the best option was to use the largest list of terms ( $LT_{99}$ ) and the KNN classifier.

It is worth noting that the categories of an article are not mutually exclusive. In fact, an article can be classified as belonging to more than one category. This might explain the greater error rates obtained in the categories *national* and *international*: one might argue that these categories correspond less to the topic covered by the article and are more related to the location of the news content.

## 5 Contributions and Future work

For copyright reasons, the corpus used to train and evaluate the classifiers described in this article cannot be shared. All the code used to process the documents, to implement the classifiers and evaluate them can be found at <http://github.com/andrefs/mapi-msr-categorization>.

The tasks and results previously described already provide useful insights into this matter. However, the lowest error rate obtained (0.35) might still be improved upon, either by leveraging new algorithms, fine tuning the ones already tested, or by increasing the train corpus.

We have established that it is in fact possible to develop a news article classifier which is lightweight enough to be used (in real time) on the client side of a web application. The following step will be the actual implementation of the classifier, probably in the form of a JavaScript library.

## Acknowledgment

André Santos has a PhD scholarship from Fundação para a Ciência e Tecnologia.

## References

1. Allan, S.: Online news: Journalism and the Internet. McGraw-Hill Education (UK) (2006)
2. Boczkowski, P.J.: Digitizing the news: Innovation in online newspapers. mit Press (2005)
3. Bollinger, L.C.: The tolerant society. Oxford University Press on Demand (1988)
4. Borko, H., Bernick, M.: Automatic document classification. Journal of the ACM (JACM) 10(2), 151–162 (1963)

5. Christopher, D.M., Prabhakar, R., Hinrich, S.: Introduction to information retrieval. *An Introduction To Information Retrieval* 151, 177 (2008)
6. Chyi, H.I., Lasorsa, D.: Access, use and preferences for online newspapers. *Newspaper Research Journal* 20(4), 2–13 (1999)
7. Goode, L.: Social news, citizen journalism and democracy. *New media & society* 11(8), 1287–1305 (2009)
8. Greer, J., Mensing, D.: The evolution of online newspapers: A longitudinal content analysis, 1997-2003. *Internet newspapers: The making of a mainstream medium* pp. 13–32 (2006)
9. Hilbert, M., López, P.: The world’s technological capacity to store, communicate, and compute information. *science* 332(6025), 60–65 (2011)
10. House, F.: *Freedom of the Press 2008: A global survey of media independence*. Rowman & Littlefield Publishers (2009)
11. Kim, S.M., Hovy, E.: Extracting opinions, opinion holders, and topics expressed in online news media text. In: *Proceedings of the Workshop on Sentiment and Subjectivity in Text*. pp. 1–8. Association for Computational Linguistics (2006)
12. O’hara, K., Sellen, A.: A comparison of reading paper and on-line documents. In: *Proceedings of the ACM SIGCHI Conference on Human factors in computing systems*. pp. 335–342. ACM (1997)
13. Rubin, T.N., Chambers, A., Smyth, P., Steyvers, M.: Statistical topic models for multi-label document classification. *Machine learning* 88(1), 157–208 (2012)
14. Sebastiani, F.: Machine learning in automated text categorization. *ACM computing surveys (CSUR)* 34(1), 1–47 (2002)
15. Teitler, B.E., Lieberman, M.D., Panozzo, D., Sankaranarayanan, J., Samet, H., Sperling, J.: Newsstand: A new view on news. In: *Proceedings of the 16th ACM SIGSPATIAL international conference on Advances in geographic information systems*. p. 18. ACM (2008)
16. Vadrevu, S., Nagarajan, S., Gelgi, F., Davulcu, H.: Automated metadata and instance extraction from news web sites. In: *Web Intelligence, 2005. Proceedings. The 2005 IEEE/WIC/ACM International Conference on*. pp. 38–41. IEEE (2005)
17. Williams, P., Leighton John, J., Rowland, I.: The personal curation of digital objects: A lifecycle approach. In: *Aslib Proceedings*. vol. 61, pp. 340–363. Emerald Group Publishing Limited (2009)
18. Yaginuma, T., Pereira, T., Baptista, A.A.: Design of metadata elements for digital news articles in the omnipaper project (2003)
19. Yaginuma, T., Pereira, T., Baptista, A.A.: Metadata elements for digital news resource description (2003)