

# Information Processing and Retrieval

Master in Informatics Engineering and Computation at FEUP, U.Porto

Vitor Cavaleiro  
up202004724@edu.fe.up.pt

Diogo Fonte  
up202004175@edu.fe.up.pt

Rodrigo Figueiredo  
up202005216@edu.fc.up.pt

Sofia Rodrigo  
up202301429@edu.fe.up.pt



Figure 1: Topic image

## ABSTRACT

The main purpose of this project is the collection, preparation and processing of a specific unstructured textual dataset, for a later development of a search system. This report intends to document this process, dividing it in three sections which distribute the preparation, retrieval and usage of the data.

## KEYWORDS

Dataset, data, news, pipeline

## 1 MILESTONE 1: DATA PREPARATION

### 1.1 Data collection

**1.1.1 Dataset choice.** Before selecting the dataset, it was essential to reach a consensus on the theme. After deliberating the pros and cons of various topics, we concluded that news would be a highly suitable subject for an information search system, given the possibility of finding multiple datasets meeting our desired quality criteria.

Once the theme was decided, we initiated the search for datasets that met the necessary requirements, particularly concerning data size and quality. The chosen dataset encompassed 204,000 articles

from 18 American publications, collected from Components<sup>1</sup>. However, in order to incorporate data from another source and introduce some complexity, we opted for an additional dataset, this time consisting of 2,555 documents sourced from the BBC website, collected from ML Resources<sup>2</sup>.

The first dataset has an MIT license, while the second one is open for academic purposes.

**1.1.2 Dataset content.** The first dataset contains 204,135 articles from 18 American publications. Includes date, title, publication, article text, publication name, year, month, and URL (for some). Articles mostly unevenly span from 2013 to early 2018, with a smattering pre-2013.

The second one consists of 2225 documents from the BBC news website corresponding to stories in five topical areas from 2004-2005, divided in 5 categories: business, entertainment, politics, sport and tech, and includes raw text data, from where we extracted the needed attributes.

To merge the datasets, we standardized both of them into the same format, ensuring the following attributes:

- **title:** Title of the article
- **author:** Author of the article

<sup>1</sup><https://components.one/datasets/all-the-news-articles-dataset>

<sup>2</sup><http://mlg.ucd.ie/datasets/bbc.html>

- After merging both datasets, we had one csv file with all the data and there were some things we needed to do in order to clean the data. The first thing was removing all duplicates, as well as replacing empty strings with NaN. We also decided to remove every row with missing title or content as that is valuable information that needs to be present. There was also a problem with the author names, as a lot of them were coming with "\n" before and after their names, so we removed that as well.

In order to extract more information about the chosen dataset and with the help of python libraries such as matplotlib, pyplot or seaborn, we have developed different data analysis plots to represent some parameters. Each of the following tries to obtain some characterization of the data to help us understand the information we are handling:



Figure 2 represents a wordcloud, which is generated from article content using keyword extraction to identify the most common words or phrases. We could conclude that our data is mostly dedicated to political issues.

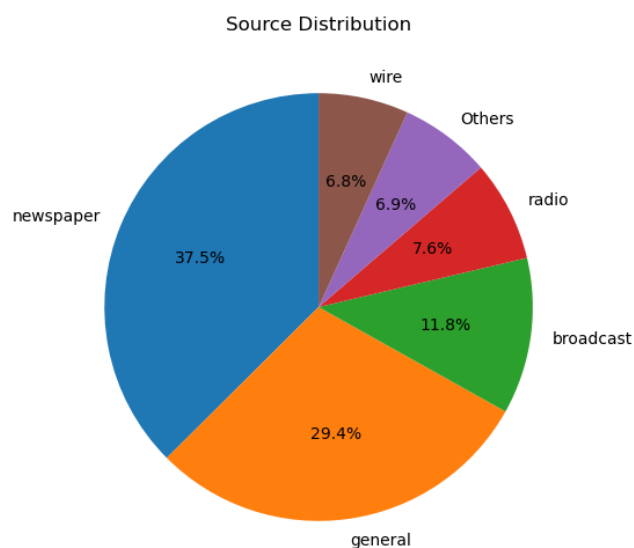


Figure 3 represents a pie chart, which visualizes the proportion of articles from each publication source within the dataset. The plot clearly emphasizes the newspaper as the most used source of news release.

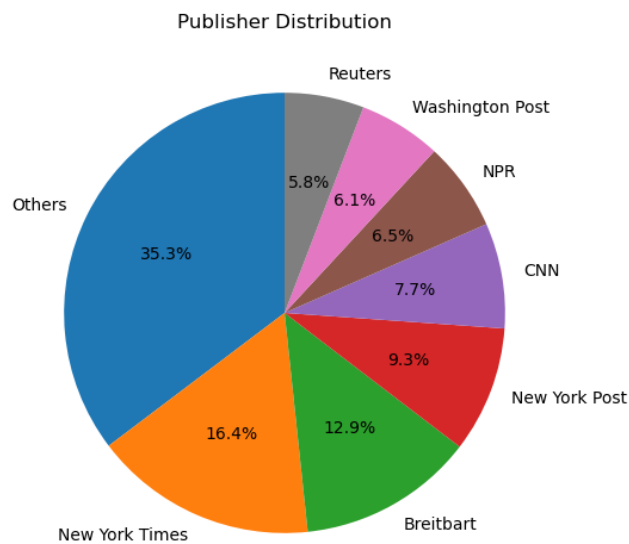


Figure 4: Publishers pie chart

Figure 4 represents another pie chart, which shows the percentage of articles from each publisher within the dataset. The most common ones are The New York Times, Breitbart and The New York Post.

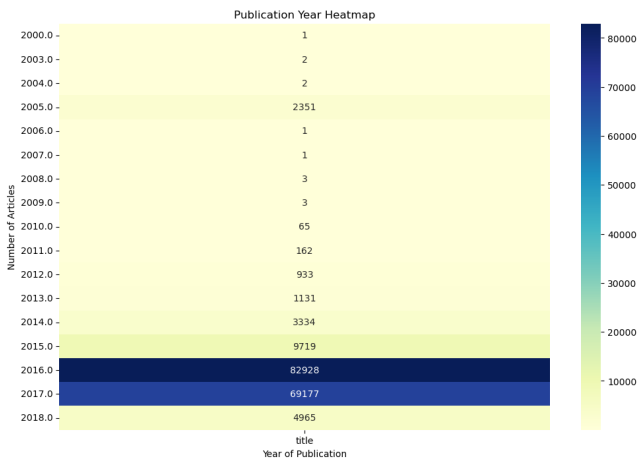


Figure 5: Heatmap of articles per year

Figure 5 represents a heatmap, a multivariable plot generated to visualize the number of articles published per year. 2016 and 2017 are clearly the years that cover the most amount of written articles.

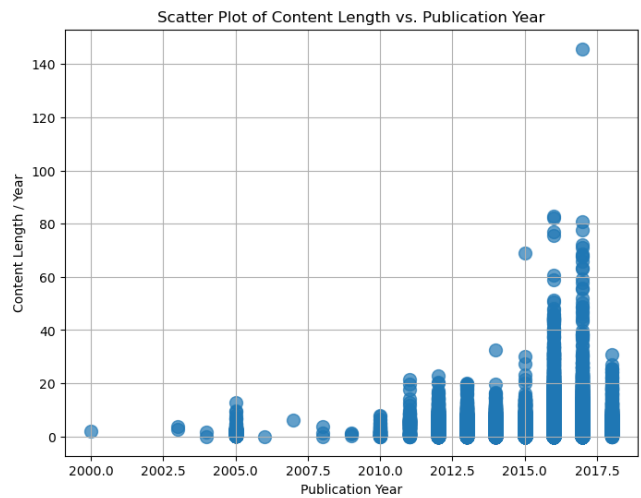


Figure 6: Scatter plot of content length per year

Figure 6 describes a scatter plot, which explores the correlation between the length of the articles and the publication date. It exposes 2015, 2016 and 2017 as the years with longer articles and also with the most amount, just as shown in the previous plot.

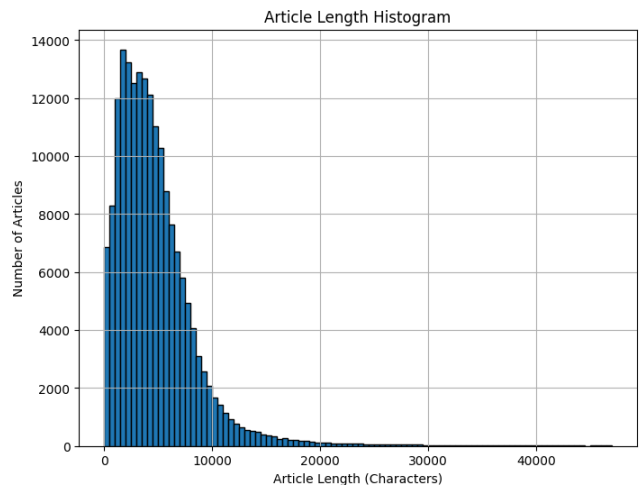


Figure 7: Article Length Histogram

Figure 7 represents a histogram of the articles length, and it can be observed that it has a right-skewed distribution with the majority of articles having a length of less than 10000 characters.

## 1.4 Domain data model

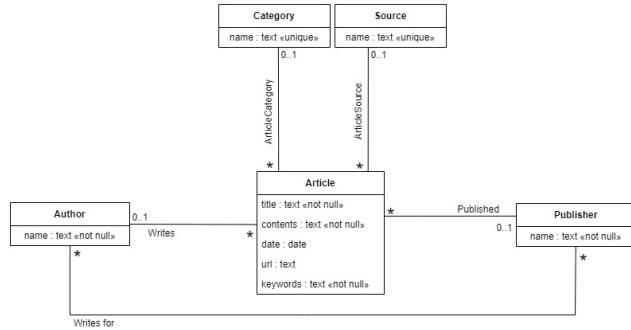


Figure 8: UML Class diagram. (<https://app.diagrams.net/>).

The data model comprises four primary classes that represent key entities within the domain: Article, Author, Publisher, and Category. These classes encapsulate essential aspects of the domain and are interrelated to facilitate effective data organization.

The Article class serves as the central entity and embodies written content within the domain. It possesses several attributes that define articles, including a title, content, publication date, URL, and keywords. Articles have associations with the Author, Publisher, and Category classes, enabling the linkage of articles to authors, publishers, and categories, as appropriate.

Authors represent individuals responsible for creating articles and the sole attribute for authors is their name. An Author class is associated with the Publisher class, illustrating the collaborative relationship between authors and publishing entities.

Publishers signify the organizations or entities responsible for publishing articles, and the only attribute is the name of the publisher. This class maintains a connection with the Author class, demonstrating the affiliation between authors and their respective publishers.

Categories are used for classifying articles based on various themes or topics, and their only attribute is also the category name. Categories are directly connected to the Article class, enabling the categorization of articles into relevant topics.

Sources are connected to the articles as they specify where a given article comes from. Their only attribute is the source name.

## 1.5 Data processing pipeline

The whole process of this first milestone is developed in Jupyter Notebook.

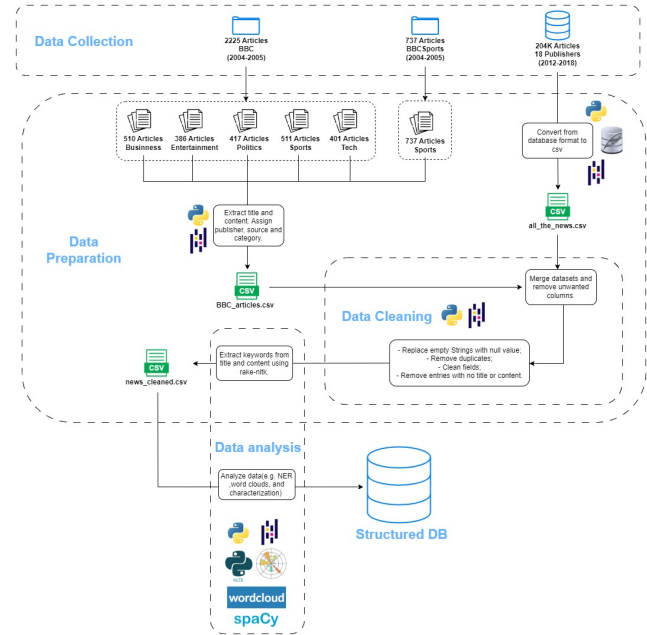


Figure 9: Data flow diagram of the pipeline. (<https://app.diagrams.net/>).

## 1.6 Prospective search tasks

The search system that is going to be developed in the next milestones requires a previous task of information needs search. These are some examples:

- Find news articles where Trump spoke on the immigration crisis
- Find news about LeBron's best performance in game
- Find articles related to latest crimes investigated by the police
- Find news articles regarding the conflicts between republicans and democrats