

News search engine

Faculty of Engineering of the University of Porto, Portugal

Vítor 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

ABSTRACT

This report intends to document the process of collection, preparation and processing of a specific unstructured news dataset, for the development of a search system. The sourced data is obtained from trustworthy sources such as Components and ML Resources.

The data preparation phase was conducted throughout a pipeline, structuring and cleaning the information from the CSV files. Subsequently, a detailed analysis was conducted, leveraging graphical representations to enhance the understanding of the dataset's inherent patterns. This first part was finalized with a statement of the prospective search tasks that were going to be worked with the following phases.

The subsequent step in this process involved indexing and filtering the information to facilitate efficient retrieval through queries in Solr. Each search task has been then thoroughly evaluated using precision-recall curves, comparing two different retrieval setups.

The final part

KEYWORDS

Dataset, data, news, pipeline, indexing, query

1 INTRODUCTION

In today's rapidly evolving information landscape, the ability to access and extract relevant information from the vast ocean of digital data has become paramount. The constant flow of information, generated by a multitude of sources, from social media platforms and news websites to academic journals and research repositories, underscores the pressing need for robust information retrieval tools.

This article introduces a cutting-edge news search engine, specifically designed to meet the challenges of information retrieval. By combining a data preparation and cleaning methodology, a further understanding of it through graphic analyzers and an indexing and retrieval tool such as Solr, this search engine promises to deliver a tailored and insightful user experience, making it an indispensable tool for individuals and researchers.

The report is divided in two principal sections.

The first one gathers the data preparation process, which includes how the datasets were selected, what they contain and the explanation of the pipeline carried out to obtain data prepared for its analyzing. This first milestone lays the groundwork for a simpler information needs retrieval.

The second section involves information indexing using Json schemas and algorithms, as well as a profound evaluation of the different proposed queries and their results.

2 DATA PREPARATION

Data preparation is the process of preparing raw data so that it is suitable for further processing and analysis. Key steps include collecting, cleaning, and labeling raw data into a form suitable for exploring and visualizing the data. Using specialized data preparation tools is important to optimize this process [4].

2.1 Data collection

Data collection is the first important step to produce a reliable search system. It involves choosing and exploring the content of the dataset, as well as exploring its quality.

2.1.1 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 [6]. 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 [5].

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

2.1.2 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
- date: Date of publication in format YYYY-MM-DD
- content: Textual content of the article
- publisher: Publisher of the article
- source: From where the article was extracted

- **category:** Describes the topic or subject matter of the article.
- **url:** Source webpage of the article

2.1.3 Data quality. The data was sourced from Components and ML Resources, two well-known and reputable platforms. It also has a large volume, is diverse and well documented.

There were some null values and missing data for some fields, but everything was handled the way it should and the data was standardized into a format we specified.

2.2 Data preparation

Data preparation is the principal phase of this first milestone. Data ingestion and data cleaning ease the management of indexing tools that will be concreted afterwards.

2.2.1 Data ingestion. It is important to clean, reduce or filter the data before its analysis and processing.

Firstly, the first dataset that comes in a .db file is exported to a JSON file using SQLite studio. Then, a python script converts through simple functions and dataframes the JSON file to a CSV file.

A similar process takes place with the BBC database. In this case, the data is collected from BBC and BBCSports articles written in various txt files, separated in 5 folders according to the categories. Afterwards, another python script combines them into a CSV file.

Finally, we bind both CSV files into one, which will be the final format of the data for further usage.

2.2.2 Data cleaning. With the help of the pandas library of python, we manipulate the CSV file throughout dataframes.

For the first dataset we start by removing irrelevant parameters such as the first column, the *id*, the *digital* and the *section* in order to transform our data into a more useful information source. Moreover, we eliminate redundant columns such as *year* and *month* since it is already detailed in the *date* parameter. Another important change was to rename some columns to more descriptive words for the actual meaning of the columns. *Publication* parameter was changed to *publisher*, *category* to *source*, and *section* to *category*.

The second dataset was exported as raw text files divided into categories, so we extracted the relevant attributes from the text and merged all the categories into one csv file, leaving it in the same format as the first dataset, so we could merge both datasets.

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.

2.3 Data analysis

At the start of this stage we decided to create a new column referring to the keyphrases of each article. We used the rakt-nltk library for that.

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 1: Articles wordcloud

Figure 1 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.

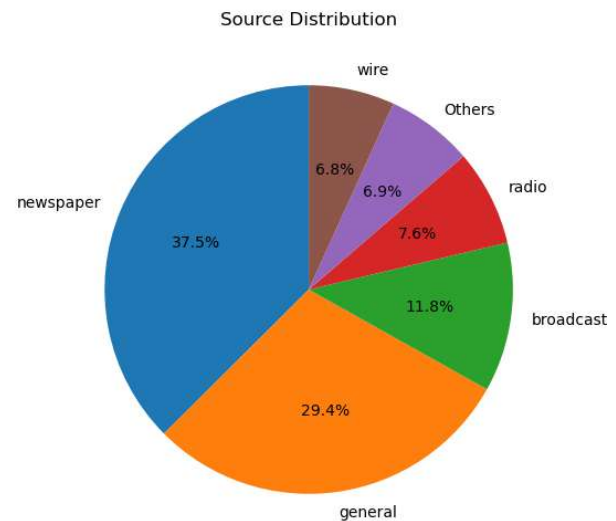


Figure 2: Sources pie chart

Figure 2 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. The general category includes articles that were published in several platforms (newspaper, website, etc).

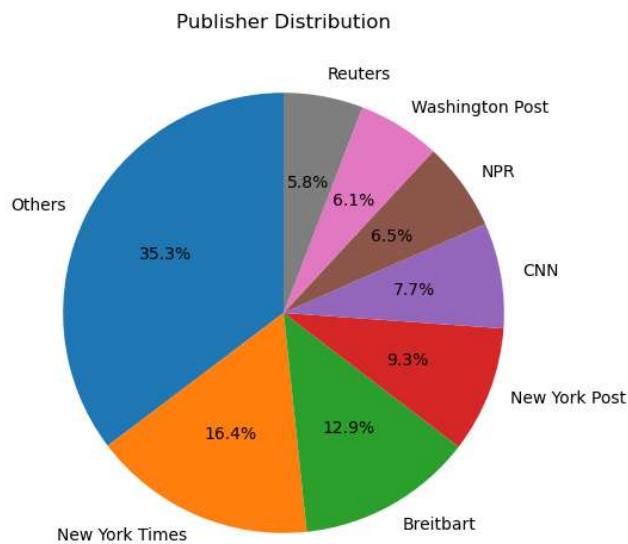


Figure 3: Publishers pie chart

Figure 3 represents another pie chart, which shows the percent of articles from each publisher within the dataset. The most common ones are The New York Times, Breitbart and The New York Post.

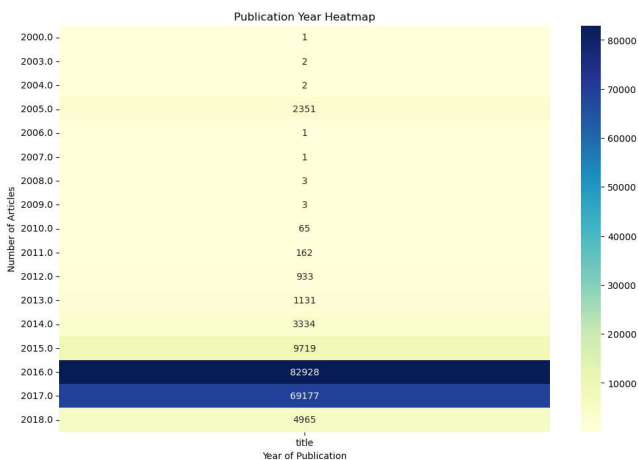


Figure 4: Heatmap of articles per year

Figure 4 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.

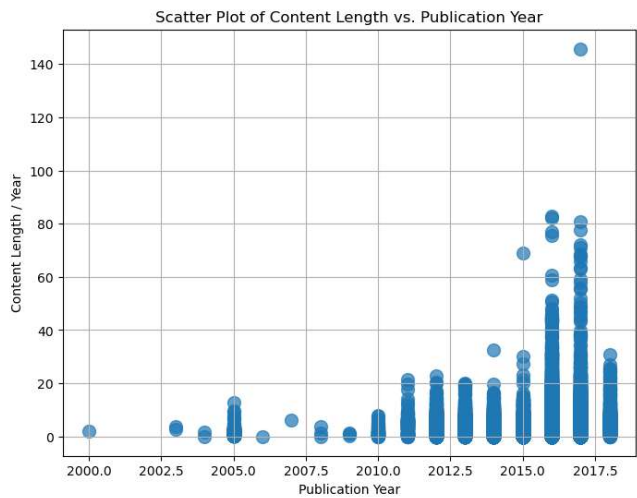


Figure 5: Scatter plot of content length per year

Figure 5 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.

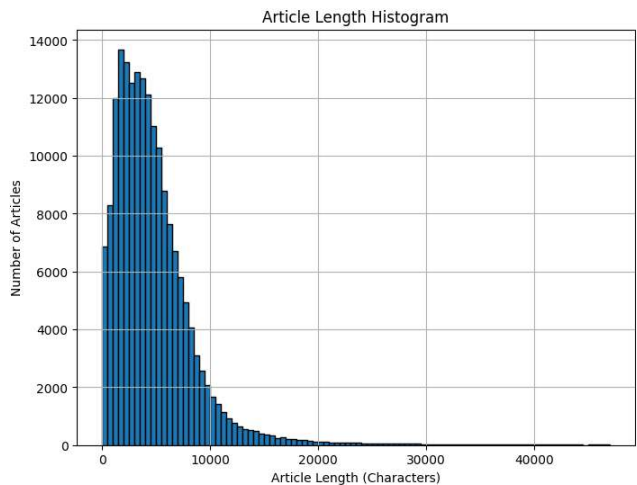


Figure 6: Article Length Histogram

Figure 6 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.

2.4 Domain data model

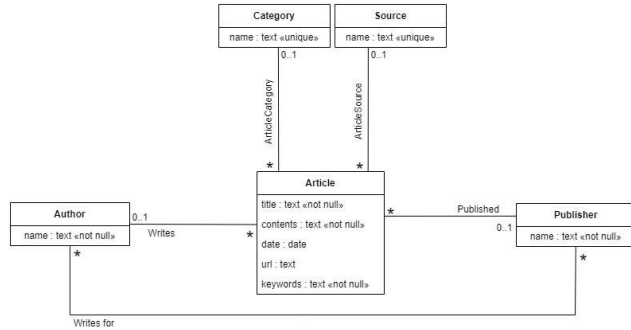


Figure 7: Domain data model [2]

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.

2.5 Data processing pipeline

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

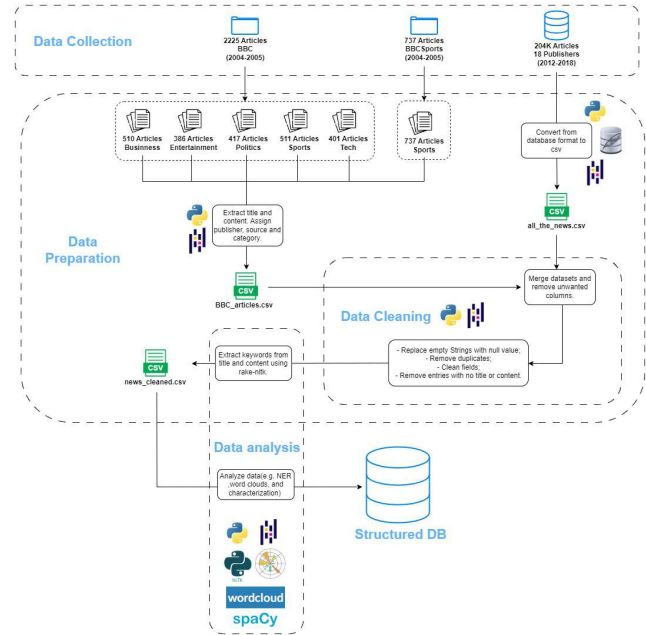


Figure 8: Data flow diagram of the pipeline [2]

2.6 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 good performances in games his team won
- Find articles related to homicides investigated by the FBI in 2017
- Find news articles regarding the conflicts between republicans and democrats about gun ownership

3 INFORMATION RETRIEVAL

Information retrieval is the activity of obtaining information system resources throughout content-based indexing, that are relevant to an information need from a collection of those resources. [7]

3.1 Tools

For this milestone, the use of Solr [1] and Docker containers [3] was essential for building our information retrieval system. This software provides better tools for indexing and querying long textual fields and also supports CSV files as input. It has also powerful search capabilities, is highly flexible and can index a wide range of document types. Solr incorporates as well advanced text analysis and tokenization features.

3.2 Indexing

To know which fields we needed to index, we focused on the domain data model and the prospective search tasks subsections. In addition, we have created a schema that contains different field types associated with each field depending on its characteristics. Table 1 describes this correspondance. The creation of a new field *Code* was necessary for the retrieval part.

Field	Field type
Code	Code id
Title	Synonym text
Author	Char text
Date	Date
Content	Content text
Publisher	Char text
Category	Synonym text

Table 1: Description of field type correspondance

Every field type has a Solr class, an index analyzer and a query analyzer. Each analyzer has different filters and a tokenizer, which in our case will be always the Standard Tokenizer from Solr.

The following items describe the several field types created and the filters that are applied to each one of them.

- **Char text**
 - ASCII Folding Filter: this filter converts alphabetic, numeric, and symbolic Unicode characters which are not in the Basic Latin Unicode block (the first 127 ASCII characters) to their ASCII equivalents, if one exists.
 - Lower Case Filter: converts any uppercase letters in a token to the equivalent lowercase token. All other characters are left unchanged.
 - Mapping Char Filter: this char filter is used for changing one string to another (for example, for normalizing á to a)
- **Synonym text**
 - ASCII Folding Filter
 - Lower Case Filter

- Synonym Graph Filter: this filter maps single or multi-token synonyms, producing a fully correct graph output.
- Porter Stem Filter: this filter applies the Porter Stemming Algorithm for English. The results are similar to using the Snowball Porter Stemmer with the language="English" argument.
- **Content text**
 - ASCII Folding Filter
 - Lower Case Filter
 - Synonym Graph Filter
 - Porter Stem Filter
 - English Possessive Filter: this filter removes singular possessives (trailing 's) from words.
 - Hyphenated Words Filter: this filter reconstructs hyphenated words that have been tokenized as two tokens because of a line break or other intervening white-space in the field text. If a token ends with a hyphen, it is joined with the following token and the hyphen is discarded.
 - Stop Filter: this filter discards or stops analysis of, tokens that are on the given stop words list. A standard stop words list is included in the Solr conf directory, named stopwords.txt, which is appropriate for typical English language text.
- **Date**
 - No filters are applied since it is a date and Solr has a class named TrieDateField for it.

3.3 Retrieval and Evaluation

In this subsection, we describe four different information needs built from the prospective search tasks. We will detail each one of them, as well as the query that we used, the boosts that were applied and the results we obtained.

Since the number of retrieved documents was high, we analyzed some of them for each query to determine which ones were more relevant. Using the given spreadsheet, we wrote their code in qrels.txt files.

We used several systems to test the information needs: firstly, we just used the indexed fields described in the previous subsection and afterwards we explored the following:

- **Fields boosts:** enhance the relevance of specific fields.
- **Term boosts:** elevate the importance of certain words.
- **Proximity searches:** emphasizes the closeness of terms.
- **Wildcards/Fuzziness:** allows for flexibility in matching variations of a term.

To evaluate the queries' results we take into account two main concepts: *Precision*, which is the accuracy of retrieved relevant documents, and *Recall*, that quantifies the proportion of relevant documents retrieved. We also consider the following evaluation metrics:

- **Average precision (AP):** value obtained for the set of documents existing after each relevant document is retrieved.
- **Precision at 10 (P@10):** value obtained from the number of recommended documents that are relevant divided by the number of recommended documents.

QUERIES

The following comparisons will be based on two different retrieval setups: the first one only uses a simple schema without filters and the second one uses our schema and several boosts that we described previously.

- **Q1:** Find news articles where Trump spoke on the immigration crisis
 - Arguments and boosts

Tag	Value
q	Trump immigration
qf	title ^{2.5} content
fl	all fields
bq	title:Trump ³ title:\\"Trump speak 2\\" 5 ^{2.5} content:Trump ^{2.5} content:\\"Trump speak 2\\" 5 ²

Table 2: Description Q1 arguments

- Precision metrics and P-R curves

Metric	Simple value	Value with boosts
AP	0.53	0.71
P@10	0.5	0.8

Table 3: Precision metrics for Q1

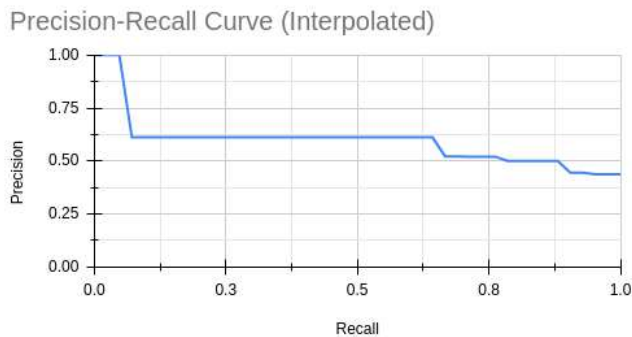


Figure 9: P-R curve for Q1 with simple schema

Precision-Recall Curve (Interpolated)



Figure 10: P-R curve for Q1 with boosts

- Discussion

As we can see in table 3, the average precision and the precision at 10 is higher for the second setup, meaning there are more relevant documents on the first positions for the second setup. Through figures 9 and 10 we can check that the second setup presents more relevant documents on the first results as the graph keeps constant at precision 1 for longer, and it also doesn't decrease that fast comparing to the first setup.

- **Q2:** Find news about LeBron's good performances in games his team won
 - Arguments and boosts

Tag	Value
q	Lebron good game win
qf	title content ^{1.5}
fl	all fields
bq	title:\\"Lebron win\\" 5 ² content:\\"Lebron bad\\" 5 ^{0.1} content:\\"Lebron fail 2\\" 5 ^{0.1} content:\\"Lebron points\\" 5 ² content:\\"Lebron win 2\\" 5 ² content:scored 2 ²
pf	title content ³

Table 4: Description Q2 arguments

- Precision metrics and P-R curves

Metric	Simple value	Value with boosts
AP	0.7	0.89
P@10	0.4	0.6

Table 5: Precision metrics for Q2

Precision-Recall Curve (Interpolated)



Figure 11: P-R curve for Q2 with simple schema

Precision-Recall Curve (Interpolated)

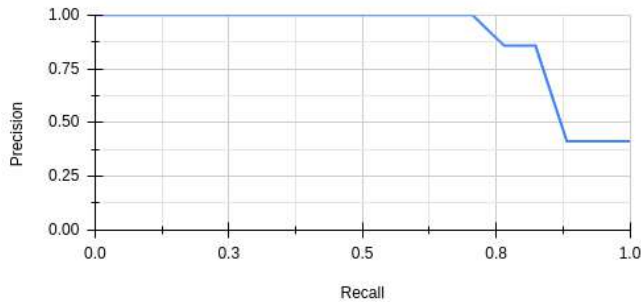


Figure 12: P-R curve for Q2 with boosts

– Discussion

As we observe in table 5 and figures 11 and 12, results are improved less than in query 1, but the second setup keeps the graph constant at precision 1 longer, meaning it has more relevant results on the top places. Precisions at 10 aren't that good but it isn't that important as we have few total relevant documents, and most of them are presented on the top.

- **Q3:** Find articles related to homicides investigated by the FBI in 2017
 - Arguments and boosts

Tag	Value
q	homicide FBI
qf	title^2 content
fl	all fields
fq	date:[2017-01-01T00:00:00Z TO 2017-12-31T23:59:59Z]
bq	title:homicides^2.0

Table 6: Description Q3 arguments

– Precision metrics and P-R curves

Metric	Simple value	Value with boosts
AP	0.62	0.81
P@5	0.2	0.6

Table 7: Precision metrics for Q3

Precision-Recall Curve (Interpolated)

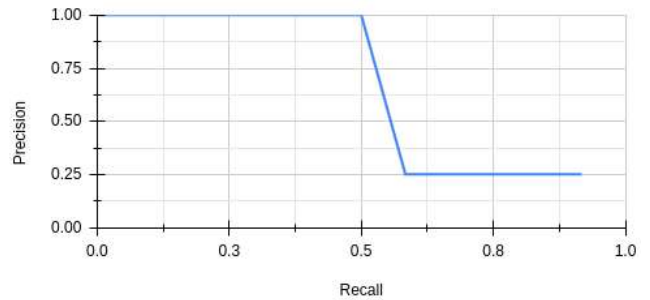


Figure 13: P-R curve for Q3 with simple schema

Precision-Recall Curve (Interpolated)

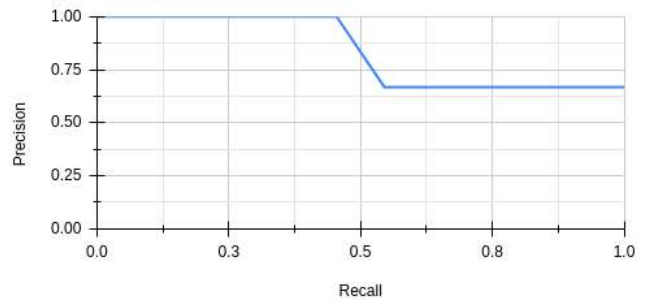


Figure 14: P-R curve for Q3 with boosts

– Discussion

For this query, as there are only 12 relevant documents in total, we chose to have the precision at 5 metric. From figures 13 and 14 we can see that the first positions have relevant documents which is a good sign, but the second setup shows more relevant documents at the last positions, that's why the graph doesn't decrease as much as the first setup. From table 7, we can check that the precisions are good, with the second setup presenting a better result for the reasons mentioned above. The precision at 5 isn't that good for the first setup as the few relevant documents presented are all gathered on the first positions, but then the following documents are mostly not relevant, making the precision go down a lot.

- **Q4:** Find news articles regarding the conflicts between republicans and democrats about gun ownership
 - Arguments and boosts

Tag	Value
q	Republicans Democrats "gun ownership"
qf	title^2.5 content
fl	all fields
bq	title:gun^1.5 content:gun^1.5 content:conflict^1.5 content:Republican 2^2 content:Democrat 2^2 content:congress^2

Table 8: Description Q4 arguments

- Precision metrics and P-R curves

Metric	Simple value	Value with boosts
AP	0.83	0.87
P@10	0.7	0.8

Table 9: Precision metrics for Q4

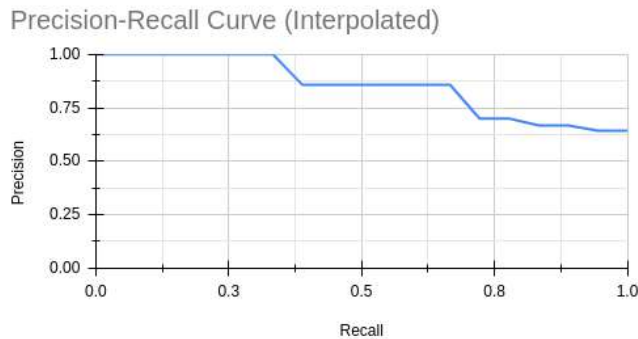


Figure 15: P-R curve for Q4 with simple schema

Precision-Recall Curve (Interpolated)

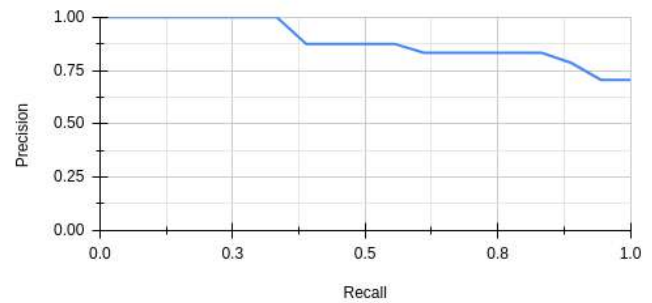


Figure 16: P-R curve for Q4 with boosts

- Discussion

As we can see in table 9 and figures 15 and 16, boosts slow down the precision decrease but do not improve final retrieval results. For this query the results were very similar as we can see in table 9, with the second setup being better by just a bit. From figures 15 and 16 we can confirm that they are very similar, presenting relevant documents on the first positions, but the second setup doesn't decrease as much, meaning the total amount of relevant documents presented is higher than the first one.

4 CONCLUSIONS

All the objectives set in the first milestone were successfully accomplished. The satisfaction of those objectives has allowed us to have a better understanding of the datasets we are working with and led to a preparation for the next stages of the project.

Regarding the second milestone, we can consider we have made possible the prospective search tasks described in the first milestone so, we can conclude that the main objective has been accomplished. Apart from this, we have to be critical of the results of the evaluation that has been carried out. The values obtained in the different queries that were executed can be improved, but still, the usage of personalized boosts has helped a lot. In every single case, the boost chosen for the query has improved the result.

REFERENCES

- [1] -. *Apache Solr Reference Guide*. <https://solr.apache.org/guide/solr/latest/index.html> (visited on 2023-11-07).
- [2] -. *Diagrams.net*. <https://app.diagrams.net/> (visited on 2023-10-10).
- [3] -. *Docker Guides*. <https://docs.docker.com/get-started/overview/> (visited on 2023-11-07).
- [4] -. *What is data preparation?* https://aws.amazon.com/what-is/data-preparation/?nc1=h_ls (visited on 2023-10-07).
- [5] 2016. *BBC dataset from ML resources*. <http://mlg.ucd.ie/datasets/bbc.html> (visited on 2023-10-09).
- [6] 2019. *News dataset from Components*. <https://components.one/datasets/all-the-news-articles-dataset> (visited on 2023-10-09).
- [7] 2020. *Information Retrieval*. <https://www.librarianshipstudies.com/2020/02/information-retrieval.html> (visited on 2023-11-13).