

Information Processing and Retrieval

Part 1 Project Report

Group 08:

Daniele Avolio, ist1111559

Michele Vitale, ist1111558

Luís Dias, ist198557

■ 1 Problem Statement

In this project we handled the task of **summarizing** and **extracing keywords** from a set of documents. In particular, we were handling documents regarding news from the BBC. Our dataset is composed both from the plain text of the news and the corresponding summarization, retrieved using state of the art techniques. Note that in this part of the project we didn't use the summarization of the news for guiding our system, but we only used it for evaluating the performance of it.

The tasks that were conducted can be explored one by one. Let's start by listing them:

- **Indexing:** The creation of a structure that allows to quickly retrieve the documents.
- **Text Summarization:** Given a document, the task is to compute the best set of sentences that are more relevant to the document itself.
- **Keyword Extraction:** Given a document, retrieve the most important words that are present in the document.
- **Evaluation:** Given a set of produced summaries S_p and a set of real summaries S_r , compute the metrics that evaluate the performance of the system.

Moreover, for the task of **Text Summarization** we explored some techniques to improve the performance of the system. Namely, we tried to use **Reciprocal Rank Fusion (RRF)**, a technique that allows to combine the results of different systems in order to improve the performance of the system itself, and **Maximal Marginal Relevance (MMR)**, a techniques that theoretically reduces the redundancy of the produced summaries.

■ 2 Adopted solutions

Each task of the project presented a set of choices or challenges that we had to overcome. This section of the report summarizes briefly those aspects, describing our solution in a non-exhaustive way. To better understand our proposed system, please refer to the code in the ipynb file.

■ 2.1 Indexing

■ 2.2 Text summarization

■ 2.3 Keyword extraction

■ 2.4 Evaluation

■ 3 Proposed questions

In this section, we answer point-to-point to the questions proposed during project description.

■ 3.1 Describe the corpus D and summaries S . Are terms uniformly distributed regarding TF-IDF?

Fixed the x axis on terms, it can be seen by the figures that the plotted distribution of words remains close to unchanged. However, it is important to note that the y axis has a different scale, due to the smaller cardinality of the summary set.

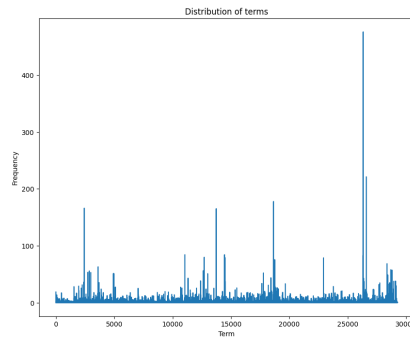


Figure 1: Corpus distribution

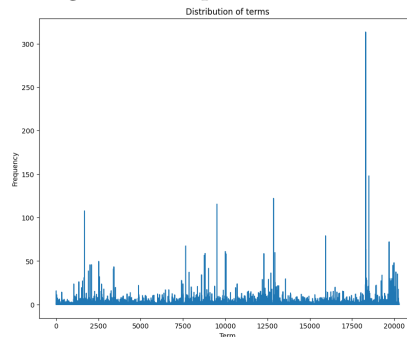


Figure 2: Summaries distribution

■ **3.2 How does the summarization system perform for the full collection? And within each category? Any intuition for the observed differences?**

QUA SI POSSONO VEDERE DEI T TEST

Non mi trovo con il numero di sample. Da controllare

The following graph shows CONCLUSIONI

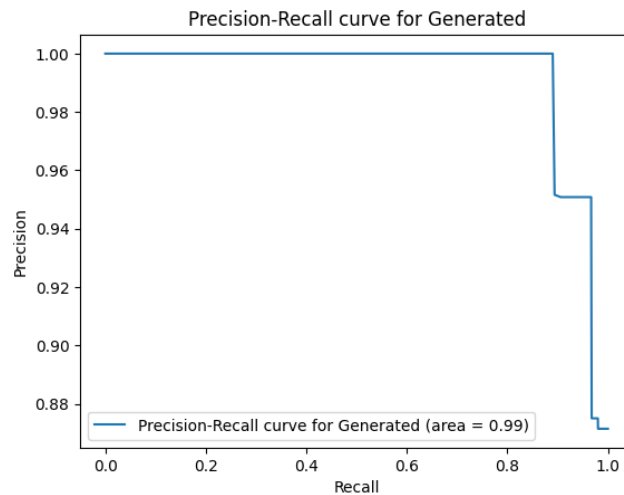


Figure 3: Performances

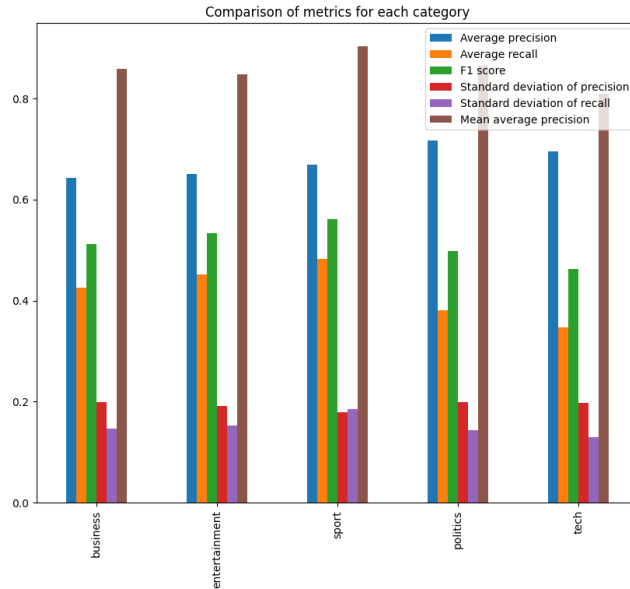


Figure 4: Performances

Category	Documents	Precision	Recall	F1	MAP
business	510	0.642	0.426	0.512	0.859
entertainment	386	0.651	0.452	0.533	0.848
sport	511	0.669	0.482	0.561	0.904
politics	417	0.717	0.381	0.498	0.864
tech	401	0.695	0.347	0.463	0.810
whole.collection	2225	0.675	0.421	0.518	0.992

Table 1: Performance metrics for different categories.

3.3 How IR models affect summaries? How vector space models compare with language models?

Qua va un attimo contestualizzata la figura

3.4 Is Reciprocal Rank Fusion (RRF) useful to aid decisions?

We expected that RRF would help to improve the summaries, since it's a way to combine the results of different systems and, in machine learning contexts,

ensembles are generally outperforming single models.

Our test was done on the entire collection due to the scarcity of computational power. Thus, we conducted a test on a subset of the dataset, that shows how the RRF has no particular difference, with the MAP being worse.

STA COSA NON HA SENSO NON HO CAPITO CHE INTENDI Probably our BERT implementation is not the best, so probably this can lead to a worse result since the BERT score is calculated in a way that I'm not sure if it's the best.

FINE COSA NOSENSO Moreover, a BERT-based solution is not suggested since the model usage times, combined with the additional overhead caused by calculating distances in the latent space, cause a significant slower response.

Method	Avg. Prec	Avg. Rec	F1 Score	Std Prec	Std Rec	MAP
BM25	0.7554	0.4250	0.5440	0.1962	0.1574	0.7869
RRF	0.5105	0.2832	0.3643	0.1223	0.1140	0.5014

Table 2: Performance metrics for BM25 and RRF.

■ **3.5 Considering MMR, how λ impacts the accuracy (against ideal extracts) of summaries? Should λ be a fixed threshold or depend on the provided topic document (d-specific)?**

Qua la figura va rivista, le lambda nell'asse x non corrispondono ai valori proposti e poi l'asse sembra su un insieme continuo. Ho provato a leggere il codice, non capisco che sta facendo

■ **3.6 At the suggested p length threshold, is the system better at promoting recall or precision?**

? che stai provando a dire?

Every time we try to measure the performance of a system, we always see that the precision is the highest, but the recall is the lowest. This is because the system is trying to avoid to put in the summary a sentence that is not relevant, but this can lead to a lower recall.

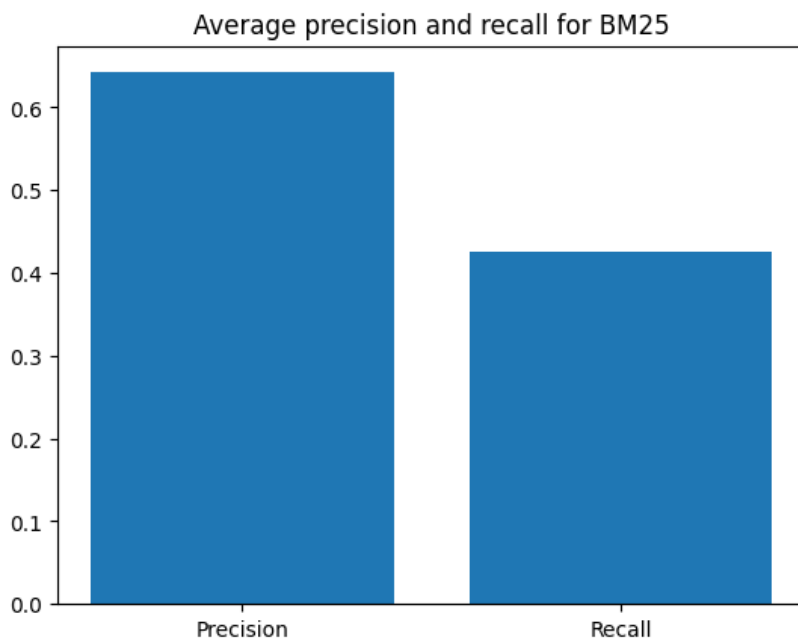


Figure 5: Precision and recall at given threshold