

Information Processing and Retrieval

Part 1 Project Report

Group 08:

Daniele Avolio, ist1111559

Michele Vitale, ist1111558

Luís Dias, ist198557

■ 1 Problem statement

In the same dataset context as Part I of this project, we are now addressing the tasks of clustering and classification with both an unsupervised approach and a supervised one, using summaries that are being provided for each document as reference.

Tasks conducted in this second part are:

- **Part A: Clustering;** for each document, the goal is grouping sentences based on their features and similarities. With this done, it is easy to select the most relevant sentences based on some criteria and algorithm that we defined.
- **Part B: Classification;** given a document, the goal is to split it into sentences and, using a binary classifier, define whether each sentence belongs to a summary or not.

This report can not contain all the data and graphs that we produced, so for more complete informations it is strongly suggested to check the comments on the provided notebook.

Some tasks were again very intensive in term of computation, so we decided to not use BERT embedding representations, since it would increase by a lot the time needed to run the code. Thus, our attention was on space representation using TF-IDF.

■ 2 Adopted solutions

Part A: Clustering

This part is conducted in an unsupervised approach, with the idea of grouping sentences with clustering algorithms. In particular, we used the **sklearn** library, with the AgglomerativeClustering class as suggested.

The main paths to explore in this part are:

- **Number of clusters and used metrics:** the main challenge here was the correct choice of the **number of clusters**, because it's a parameter that could have a big impact on the results. Using **silhouette score** we have solved this problem in an iterative way.
- **Sentences selection:** the second challenge was to select the sentences that would be used to build the summary. Using **centroids** of each cluster, we were able to build summaries that had more topics and were more representative of the original text.

Number of clusters and used metrics

A problem related to part A is to define a good number of clusters to represent the feature space of the document.

In this part, we tried to construct a custom metric taking into account Silhouette score, Calinski-Harabasz score and a function of the number of clusters. The idea was to consider the number of clusters, that can vary in the range $[2, numberOfSentences]$, to avoid high sparse clusters representation with single-sentence clusters, but in the end we noticed that the silhouette score by itself was performing overall better.

We are leaving the code for the custom metric in the notebook, but it is commented since it was not used in the final version of the code.

Sentences selection

Another relevant problem, once completed the clustering part of the project, was to decide the criteria to pick the sentences to use in the summarization. We decided to use the **centroid** of each cluster and, based on the distance from it, we pick the required number of sentences. We have done some tests on different criteria, but others ideas that we had were not really convincing on summaries, so we kept the centroid distance as metric.

It is important to note that this strategy has a limitation, since with high numbers of sentences picked from each cluster there might be some redundancy in summaries. Nevertheless, with a low number of picked sentences it is leading to convincing results, taking into the summary relevant sentences.

Part B: Classification

In this section the task is to create a binary classifier that has as goal to discern if a sentence should belong to the summary or not. We used models from the **scikit-learn** library, and in particular the machine learning techniques that we exploded are Random Forest, Gradient Boosting, Gaussian Naive Bayes, K-Neighbors and Multi-layer Perceptron.

Before the training phase, the challenge was to find the best features and the correct shape to represent data and build the models. We ended up with the following structure.

similarity	n_sentence	n_words	n_stopwords	n_keywords	length_of_sentence	tfidf_score
position_in_doc	category	n_nouns	n_verbs	n_adjectives	n_adverbs	id
summary(target)						

We then analyzed the features with a correlation matrix, the Shapley value of each features via the **SHAP** library and the Scikit-learn built-in feature importance. This process lead us to drop some features, with the following schema being the one used to train our models.

similarity	n_sentence	n_words	n_stopwords	n_keywords	length_of_sentence
tfidf_score	position_in_doc	n_nouns	n_verbs	n_adjectives	n_adverbs
category_business	category_entertainment	category_politics	category_sport	category_tech	Summary(target)

At this point, we could train our models and evaluate them with the common machine learning evaluation metrics.

■ 3 Proposed questions

Part A: Clustering

■ 3.1 Question 1

Do clustering-guided summarization alters the behavior and efficacy of the IR system?

To answer this question we ran the **clustering based** algorithm using the same set of documents used in the *first part of the project*. The result shows a pretty big difference between the two approaches.

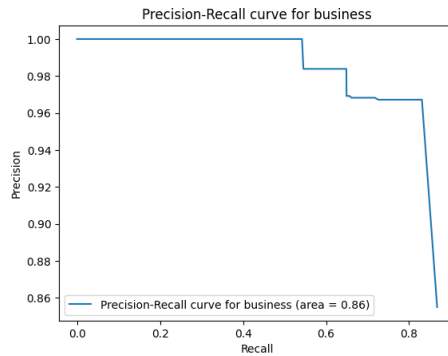


Figure 1: First Part Approach

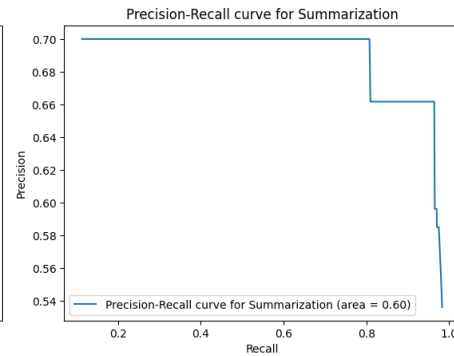


Figure 2: Clustered Approach

This result doesn't guarantee that the **clustering based** approach is worse than the approach using TFIDF and BM25, because this is based on our personal implementation of the algorithm. However, it is clear that the clustering approach is not as effective as the first approach in this case. A possible way to improve the algorithm could be to consider other sentences choices rather than focusing on the distance from the centroid of each cluster.

■ 3.2 Question 2

How sentence representations, clustering choices, and rank criteria impact summarization?

We benchmarked the performance of the clustering algorithm using a set of metrics:

Max Clusters	Number of Sentences	Our Metrics
2	3	cosine euclidean
3	5	
4	7	
6	9	
8	11	
10	13	

We didn't include any **different representations** because we only used the **TFIDF** representation. The result are indicating a very low performance of the algorithm using a specific set of metrics.

Table 1: Results of the clustering algorithm using different metrics

#clusters	#sentences	metric	avg_prec	avg_rec	f1	m_a_p
2	3	cosine	0.453504	0.449012	0.451247	0.646069
2	3	euclidean	0.453504	0.449012	0.451247	0.646069
2	5	cosine	0.457060	0.633170	0.530891	0.561392
2	5	euclidean	0.457060	0.633170	0.530891	0.561392
2	7	cosine	0.469898	0.768889	0.583312	0.492352
...
10	9	euclidean	0.484872	0.934921	0.638568	0.418472
10	11	cosine	0.486431	0.941682	0.641495	0.344811
10	11	euclidean	0.486431	0.941682	0.641495	0.344811
10	13	cosine	0.486635	0.943571	0.642110	0.336604
10	13	euclidean	0.486635	0.943571	0.642110	0.336604

More insight on this can be found in the *notebook* file.

■ 3.3 Question 3

Are anchor sentences (capturing multiple topics) included? And less relevant outlier sentences excluded? Justify

Since our algorithm is based on the **distance from the centroid** of each cluster to select the sentences, we are not able to handle the **anchor sentences** and the **outlier sentences**. We are not able to give a clear answer to this question, but a possible way to handle this could be to consider the **distance from the centroid** and the **distance from the other sentences** inside other

clusters. Sentences that are **more far** from the centroid of the cluster could be very **relevant** and could be considered as **anchor sentences**, since that sentence could be holding information between more topics.

■ 3.4 Question 4

Given a set of documents, plot the distribution of the number of keywords per document. Are keywords generally dissimilar? If not, how would you tackle this challenge? For this question we decided to use documents from

500,700

as range. The result shows that the distribution of the number of keywords per document is not uniform.

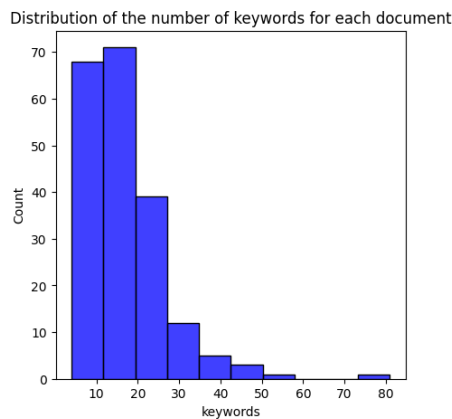


Figure 3: Distribution of the number of keywords per document

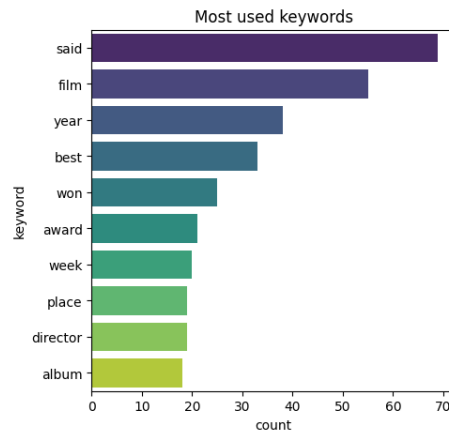


Figure 4: Distribution of the number of keywords per document

Part B: Supervised IR

■ 3.5 Question 1

Does the incorporation of relevance feedback from ideal extracts significantly impact the performance of the IR system? Hypothesize why is that so.

■ 3.6 Question 2

Are the learned models able to generalize from one category to another? Justify.

■ 3.7 Question 3

Which features appear to be more relevant to the target summarization task? Do sentence- location features aid summarization?

■ 3.8 Question 4

In alternative to the given reference extracts, consider the presence of manual abstractive summaries, can supervised IR be used to explore such feedback? Justify