

Information Processing and Retrieval

Part 2 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 $[2, \text{numberOfSentences}]$, 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.

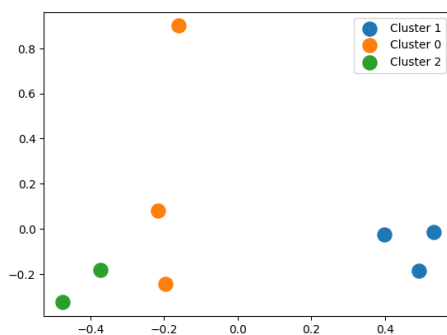


Figure 1: Example of clustering

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.

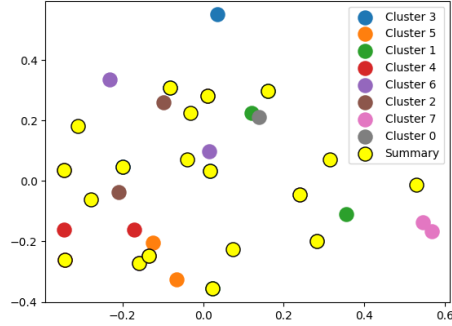


Figure 2: Example of sentences selection

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 explored 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. Since the classifier should classify sentences, each row of our dataset represents one sentence in the original library. 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)						

Note: the summary(target) contains as value [0,1] where 1 means that the sentence should be in the summary and 0 means that it should not.

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)

Correlation matrix and feature importance extraction, with related graphs, are only addressed on the notebook to save space for more important steps of the project. The main reason of features drop is high correlation shown by the analysis.

At this point, we could train our models and evaluate them with the common machine learning evaluation metrics. The one performing better is, in our case, RandomForest.

More on this can be found in the notebook, as well as in the following chapter.

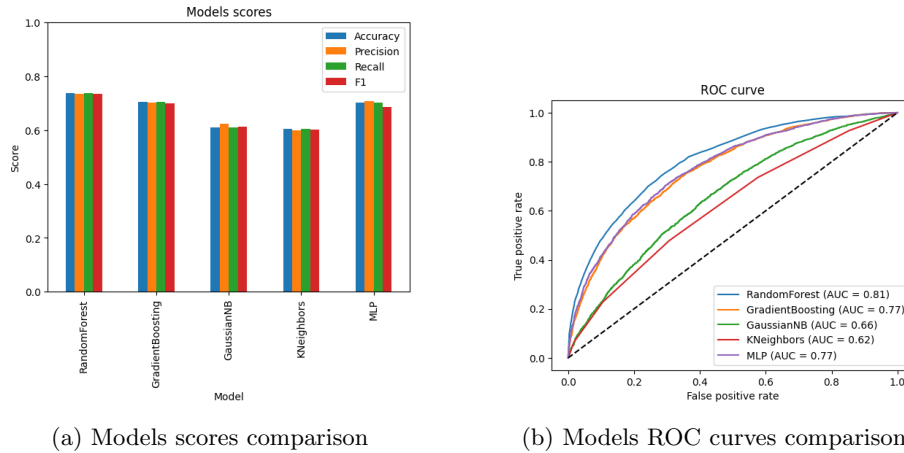


Figure 3: Scores and ROC curves

Note: The evaluation of the model is discussed in the questions.

■ 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 significative difference between the two approaches.

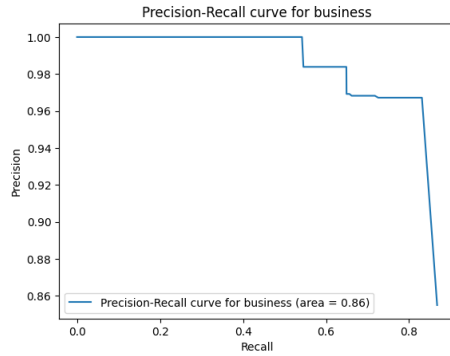


Figure 4: First Part Approach

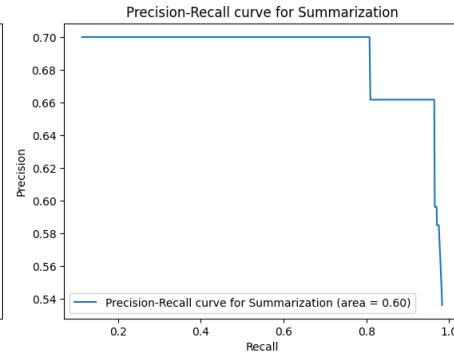


Figure 5: Clustered Approach

This result is not enough to 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 selection criteria 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	1	Cosine sim. (single) Cosine sim. (complete)
3	2	
4	3	
6	4	

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

clusters	num_sent	metric	avg_prec	avg_rec	f1	std_prec	std_rec	m_a_p
2	1	single	0.403860	0.153016	0.221942	0.171560	0.099392	0.484703
2	1	complete	0.427703	0.162918	0.235957	0.170426	0.096330	0.517471
2	2	single	0.403429	0.230112	0.293063	0.159069	0.134892	0.686571
...
4	3	single	0.405465	0.329316	0.363445	0.133851	0.193184	0.666388
4	3	complete	0.447709	0.503877	0.474135	0.110074	0.217150	0.632657
4	4	single	0.413300	0.397615	0.405306	0.119253	0.205291	0.670563
4	4	complete	0.448957	0.590894	0.510239	0.095870	0.223584	0.554584

Table 1: Results of the clustering algorithm using different metrics

The complete table can be seen on the notebook. However, we can notice that evaluation scores changes a lot in function of these parameters. In particular, the best setup for this subset of documents is obtained at max_clusters set to 6 and number of sentences per cluster set to 2, with single metric.

■ 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**. Thus, we can not give a clear answer to this question, but a possible way to deal with this problem this could be to consider the **distance from the centroid** and the **distance from the other sentences** inside other clusters. Sentences that are **farther** from the centroid of the cluster could be very **relevant** and could be considered as **anchor sentences**, since that sentence could be holding cross-cluster informations.

■ 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** to **700** as range. The result shows that the distribution of the number of keywords per document is not uniform.

Distribution of the number of keywords for each document

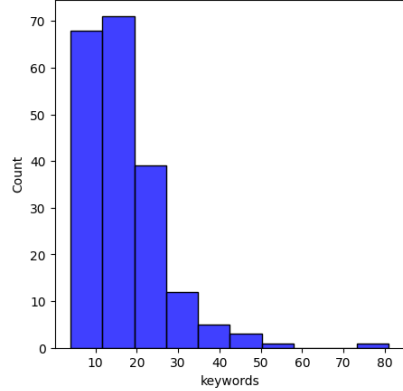


Figure 6: Distribution of the number of keywords per document

Most used keywords

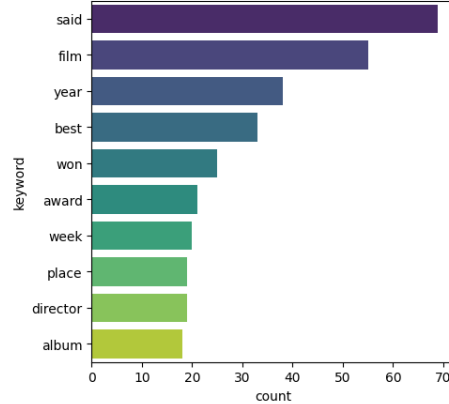


Figure 7: Top 10 used keywords

The main problem is that the **keywords** are not **dissimilar** and this could be a problem for the **clustering algorithm**. In fact, we can see that a lot of keywords are repeated in the documents. A possible way to tackle this challenge could be to use a different **keyword extraction** method. Moreover, using multiple documents of the same category could help to improve the performance of the algorithm since the keywords are more likely to be repeated in the same category.

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.** Comparing the evaluation of two best models, respectively the unsupervised model from question 3.2 and the supervised RandomForest model, it seems that the supervised IR is better performing in a perceptible way. We can say that even if the actual test sets are different since they have a significant high number of documents and the scores are way too different.

Model	Avg Precision	Avg Recall	F1 Score	#docs
Clustering	0.410742	0.304867	0.349973	200
Random Forest	0.746500	0.790460	0.767851	443

Table 2: Performance Metrics of unsupervised best model and supervised best model

From our experience we carefully can say that overall an approach based on a reference target category is usually better performing than an unsupervised model. This obviously stands assuming the high quality of the reference summaries.

■ 3.6 Question 2

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

To answer this question, we tried to train the model using the **tech category** as a training set, and all the other categories as a test set. We used *Random Forest* as a model, since it was the best between the others. The results show that the model achieved ≈ 0.69 as its best performance and ≈ 0.65 as its worst performance. This result is not bad, but it's not good either. This happens because each category is different from another, and the set of features that we are using may not be the best. However, to *effectively test the model*, what we could do is to train another model using a training set of the same size as the *tech category training set*, but using a mixture between all the categories. In this way we could have a better vision of the model's performance in the generalization task. The code and the specific result for each category can be found in the *notebook* file.

■ 3.7 Question 3

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

The most important feature that we found is the score given by the *similarity* between the sentence and the document. To get this insight, we used the *SHAP* library that uses, to highlight the importance of each feature, the concept of *Shapley Values*, a metric from game theory that given is able to assigns to each agent (our feature) a value that represents the contribution to the model's prediction. From this analysis, we can say that the *sentence-location* feature contributes to the task, since it is the 4th most important feature in the model.

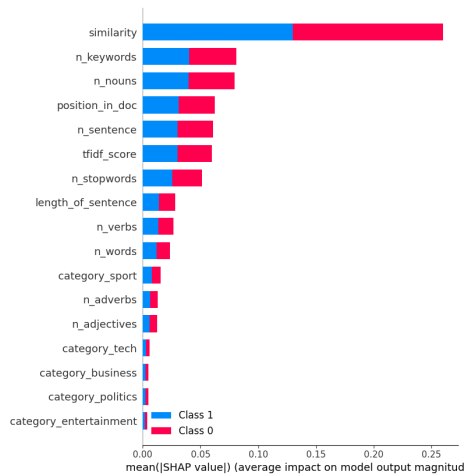


Figure 8: SHAP values for the Random Forest model

■ 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. This task might result difficult to be conducted with a machine learning approach, since there is no direct access to a target function for sentences, neither a way to build up a summary from a set of features without human supervision. This second fact makes the only possible approach to summarization being the extractive one.

The task of supervised summarization can not be conducted as we did during the part B of this delivery, since we can not construct a dataset that use sentences as features.

The abstractive approach is a summarization based on the context of the document, trying to reproduce the same concept in a different shape. This ends up in sentences from the text not being in the summary and, thus, the dataset we have created can not be replicated in the exact same shape.

A possible solution to the task with abstractive summaries could be to use the deep learning BERT, a bidirectional encoder based on Transformer architecture that is able to transpose the input text into a latent space, in which similar concepts are close to each other.

Exploiting this, our suggestion is to use the BERT model to encode both the summary, that serves as an anchor reference as provided ground truth, and the sentence we want to classify. The resulting embeddings can be used in a classifier or to compute other metrics, like distance ones, and then use those features to select sentences based on a criterion. This will still produce an extractive summary.

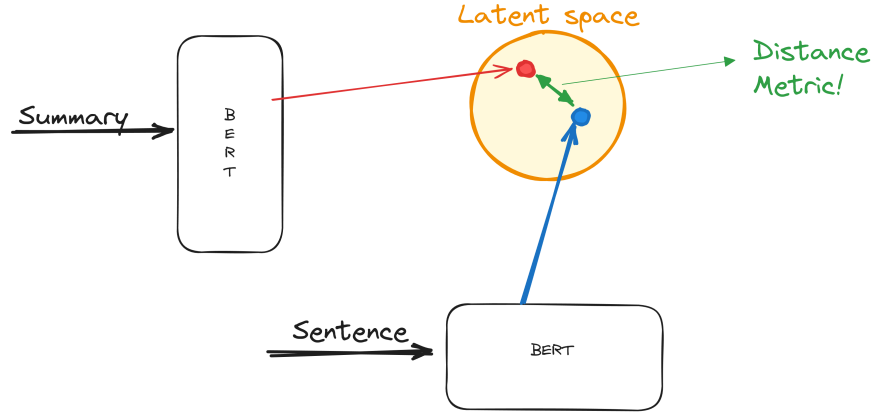


Figure 9: Possible architecture using abstractive summaries

Another approach that can be explored is based on keywords from the summary. This can be divided in two steps: during the first one, we construct a reference set of keywords, extracting them from the whole set of summaries and expanding them with one iteration on a predefined thesaurus, to make sure to have also synonyms of relevant words.

In the second phase, a sentence can be classified calculating a score which is the rate of number of keywords from the sentence that are present in the summary-extracted ones (also referred as **keywords batch** from now on) over the number of keywords in the document that contains the sentence, as it follows.

$$score = \frac{|keywordsFromSentenceInBatch|}{|keywordsInDocument|}$$

Note that this formula has the number of keywords in the document as denominator to prefer, among the same document, sentences with a higher keywords "density".

Once calculated the score, it is possible to use it to select the sentences that will be part of the summary using a fixed threshold, that can eventually be updated based on a feedback system that collects them from the final user. Here it is presented a pseudo code to better explain the algorithm.

```

1 keywords_batch = get_keywords_from_summaries()
2
3 for sentence in document:
4     score = calculate_score(sentence, document, keywords_batch)
5     if score >= threshold:
6         summary.add(sentence)
7
8 return summary

```