

# Deliverable Lab 3 CAIM

Arnau Cinca Roca, Enric Rubio Pacho

November 7, 2019

## 1 Rocchio's rule

### 1.1 Introduction

The goal of this session is to program a script `Rocchio.py` which implements a User Relevance Feedback. In order to achieve our goal we are going to make use of the Rocchio's rule.

Rocchio feedback approach was developed using the Vector Space Model. The algorithm is based on the assumption that most users have a general conception of which documents should be denoted as relevant or non-relevant. Therefore, the user's search query is revised to include an arbitrary percentage of relevant and non-relevant documents as a means of increasing the search engine's recall, and possibly the precision as well. The number of relevant and non-relevant documents allowed to enter a query is dictated by the weights of the  $\alpha$ ,  $\beta$ ,  $\gamma$  variables listed below in the Algorithm section.<sup>1</sup>

However, since we are not going to ask the user which of the queried documents are relevant or not, our script will implement Pseudo-relevance Feedback. So we have to compute the following equation:

$$Query' = \alpha \cdot Query + \beta \cdot \frac{\sum_{i=1}^k d_i}{k} \quad (1)$$

Where the second part of the equation is computed by adding the tfidf vectors from each file obtained from the *Query*.

### 1.2 Experiments

In order to implement the script, we make an *nrounds* loop. In each iteration we perform a query to *elastic search* and with the given files, we compute the tfidf vector using functions programmed in the previous session. However, to add those vectors we used dictionaries since merging vectors will have a cost of  $\mathcal{O}(n \cdot \log(m))$  instead of  $\mathcal{O}(n)$ . Furthermore, to create the *Query'*, we get the  $R$  higher weights from the  $k$  most relevant documents and finally we build the equation (1).

Also, we created a dictionary that stores the *Query* (words with its weights).

### 1.3 Observations and Conclusions

Once implemented Rocchio's rule, we ran our script with the newsgroup's collection and with *toronto* as query. The results obtained went from 359 documents on the first round to 7 on the second

---

<sup>1</sup>Rocchio's rule definition extracted from [https://en.wikipedia.org/wiki/Rocchio\\_algorithm](https://en.wikipedia.org/wiki/Rocchio_algorithm)

and, finally, we got the same number of documents in the third query as in the second one, which we can clearly see that using Rocchio's rule improves our precision in the searched query.

Furthermore, we executed it with *toronto* and *science* and found that the obtained documents talk about polithics, but there the word *science* only occurs in *Department of Computer Science of University of Toronto*. This is because we are not asking for the user opinion. Moreover, the resulting documents are contain the same text so the most relevant words obtained have three times more relevance.

Finally, we try to find the optimal parameters ( $nrounds$ ,  $R$ ,  $k$ ,  $\alpha$  and  $\beta$ ), experimenting with this parameters, we observe that in the second round ( $nrounds=3$ ), the queries converges, someones even in the first round. We also test  $R$  values and we conclude that higher values tend to return zero documents, and with smaller values tend to return more documents. With *toronto* as query, with a  $R$  of 3-4 the number of returned documents are seven, and with a higher value returns zero documents. Modifying the value of  $k$  we found that taking more documents makes the query more general and more related with the initial query, with a smaller  $k$ , the result is more focused in the theme of the top documents of the previous query, this might result in a unprecise search, because the user may not be interested in this theme. Testing  $\alpha$  and  $\beta$ , doesn't change the query, we suppose that is because in the first round the query reduces the results draslicly and this makes that  $\alpha$  and  $\beta$  only can change the relevance of the documents inside the *Query*', with a bigger data,  $\alpha$  and  $\beta$  may be usefull to improve the result.