

In this session:

- We will program a simple User Relevance Feedback cycle on top of Elasticsearch.
- We will evaluate this strategy over a collection of documents.

1 Document relevance

To show more precisely how Elasticsearch works, you have a script named `SearchIndexWeights.py` that it is similar to the script from the first session, but now only shows a specific number of hits (it shows the total number of documents that matches the query as the last line of the output) and also has a relevance score computed by Elasticsearch representing how well the document matches the query. The script has a `--nhits` flag that changes how many hits are shown and a `--query` flag that accepts a list of words (performs an AND query with all the words), this flag has to be the last one when invoking the script.

Play a little bit with this script, performing different queries and observing the scores. The syntax of the query allows using the fuzzy operator `~n`, but also the boost operator `^n` with n indicating how important this term is relative to others, which has an effect on the resulting relevance score. For example, with the `20_newsgroups` corpus try the following queries:

```
python SearchIndexWeights.py --index news --nhits 5 --query toronto nyc
python SearchIndexWeights.py --index news --nhits 5 --query toronto~2 nyc
python SearchIndexWeights.py --index news --nhits 5 --query toronto nyc^2
```

You will see that the scores and the positions of the documents change. Invent new queries and observe the results.

2 We will, we will Rocchio you

The goal of this session is to program a script `Rocchio.py` that implements a User Relevance Feedback system using Rocchio's rule. In fact, we will implement Pseudo-relevance Feedback since we will not ask the user which documents she finds relevant: simply we will assume that the first k documents are the relevant ones, for a k of our choice.

More precisely, we want our script to do the following:

1. Ask for a set of words to use as query
2. For a number of times (*nrounds*):
 - (a) Obtain the k more relevant documents
 - (b) Compute a new query applying Rocchio's rule to the current query and the Tf-Idf representation of the k documents
3. Return the k most relevant documents after the n iterations

You will have to implement the function that computes the new query applying Rocchio's rule, namely¹:

$$Query' = \alpha \times Query + \beta \times \frac{d_1 + d_2 + \dots + d_k}{k}$$

You will have to compute the tf-idf vector for each document. To average all the documents you will have to sum vectors with many elements, you must use dictionaries to do it more efficiently instead of merging ordered vectors. What is the difference in computational cost of merging ordered vectors or using dictionaries for this operation? Discuss it in your report.

Also the resulting list of terms will be large. Consider pruning the list to only the R more relevant terms (larger weights).

Most of the elements to solve this you already have from the provided scripts and from your solutions from the previous sessions:

1. From `SearchIndexWeights.py` you have the code for building a query for a list of words and then retrieving the k more relevant documents.
2. Adding the weights computed using Rocchio's rule to each word in the search needs only concatenating the word, the boost operator (\sim) and the weight.
3. You will have to compute the tf-idf vector for a document, problem that you have solved in the previous session.

Observe that there are several parameters you can play with, at least:

- $nrounds$, the number of applications of Rocchio's rule
- k , the number of top documents considered relevant and used for applying Rocchio at each round
- R , the maximum number of new terms to be kept in the new query
- α and β , the weights in the Rocchio rule.

Please make sure that they are easy to change in the code (e.g., their values defined only once!), so that you can run your experiments easily.

3 Experimenting

Once you are done with your programming, try it out with the test collections from the previous sessions. Do the queries that pseudorelevance feedback produce make sense? For example, do the new terms seem related to what the user is looking for?

In which sense do the results improve? Recall? Precision?

Investigate to some extent the effect of each parameter. Do you get very different results if you change the parameters $nrounds$, k , R , α , β ? Do you find, for each one, some value or value range that seems to be optimal in some sense?

¹We use a simplified version of Rocchio's rule setting $\gamma = 0$.

4 Rules of delivery

1. No plagiarism; don't discuss your work with other teams. You can ask for help to others for simple things, such as recalling a python instruction or module, but nothing too specific to the session.
2. If you feel you are spending much more time than the rest of the classmates, ask us for help. Questions can be asked either in person or by email, and you'll never be penalized by asking questions, no matter how stupid they look in retrospect.
3. Write a short report with your results and thoughts. Make it at most 2 pages. Strive to summarize what new things you learned in this session. You are welcome to add conclusions and findings that depart from what we asked you to do.
4. Turn the report to PDF. Make sure it has your names, date, and title. Create a single .zip file all the python scripts that you created or modified; in the modified scripts, make sure you mark visibly with comments the parts that you modified.
5. Submit your work through the Racó. There will be a Practica open for each report. Both members of the team must submit identical PDF files.

Deadline: Work must be delivered within 2 weeks from the end of the lab session. Late submissions risk being penalized or not accepted at all. If you anticipate problems with the deadline, tell us as soon as possible.