

Document retrieval task on controversial topic with Re-Ranking approach

Notebook for the Touché Lab on Argument Retrieval at CLEF 2021

Andrea Cassetta^a, Alberto Piva^a and Enrico Vicentini^a

^aUniversity of Padua, Italy

Abstract

This paper is the report of the work done for Argument Retrieval CLEF 2021 Touché Task 1 by Shanks team (based in Italy and precisely the members are University's of Padua students). Argument Retrieval CLEF 2021 Touché Task 1 focuses on the problem of retrieving relevant arguments for a given controversial topic, from a focused crawl of online debate portals. After some tests Shanks group has decided to parse the input documents taking only the title, premises and conclusion of the arguments (as well as the stance necessary to understand the arguments' author point of view). After the indexing part of the documents the work is concentrate on how the retrieving and the raking are done. After some tests, we discover that the better results are obtained using a WordNet [1] based query expansion approach and a re-ranking process with two different similarity functions. This report describes in details how the documents parsing work and how the indexing and searching part are developed. The unexpected update of the qrels file did not allow us to re-run all the tests. In the end, however, we also reported the results of the runs obtained from parameter tuning on the new qrels.

Keywords

Argument Retrieval CLEF 2021 Touché Task 1, WordNet synonyms, Re-ranking, BM25, DirichletLM

1. Introduction

In this report, we describe the project developed for the participation by the Shanks group to the CLEF 2021 Touché Task 1. The task focuses on the problem of retrieving relevant arguments for a given controversial topic, from a focused crawl of online debate portals.

Our goal is to develop a Java based information retrieval system that finds and ranks the relevant documents from the args.me corpus dataset composed by over 380.000 arguments crawled from 5 different debate forums [2] for 50 topics (query). The retrieved results need to be relevant for each input topic the system has to elaborate.

This paper is structured as follow: Section 3 is about the solutions that we've taken in consideration to build the retrieval system, Section 4 describes the whole workflow of the program;


CLEF 2021 – Conference and Labs of the Evaluation Forum, September 21–24, 2021, Bucharest, Romania
“Search Engines”, course at the master degree in “Computer Engineering”, Department of Information Engineering,
University of Padua, Italy. Academic Year 2020/2021

✉ andrea.cassetta@studenti.unipd.it (A. Cassetta); alberto.piva.8@studenti.unipd.it - 0000-0003-0242-0749
(ORCID) (A. Piva); enrico.vicentini.1@studenti.unipd.it (E. Vicentini)



© 2021 Copyright for this paper by its authors.

Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

Section 5 explains our experimental setup including the software, tools and methods used; Section 6 discusses results; and finally, Section 9 draws some conclusions and outlooks for future work.

2. Related Works

To create our search engine we build upon some source code created by Professor Nicola Ferro to show as some toy examples and changed as described in the subsequent sections. We have also read the overview of CLEF 2020 on the Touchè task[3].

After some research and as suggested by Professor Ferro, we discover a different way of re-ranking and merging the results. In practice the paper *Combination of Multiple Searches* by Fox and Shaw [4]. Shaw, described an interesting method to increase the performance of our system by combining the similarity values from different output runs, using Boolean retrieval methods. Into the paper is also described how they have done the indexing (and analyzing) part, but we decide to overtake that part because the stating dataset is different, and we have already analyzed our dataset to reach better index possible. For the purpose to understand how the results merging has been done is not useful to analyze also how the query are written and so we decide only to relate that the P-norm queries are written using a complex boolean expression using AND and OR operators. When all the runs are done, the second part of the experiment consists in combining the output runs (obviously obtained from the same collection of data) to reach the best result. Different way of combining them are, for example, taking the top N documents retrieved for each run or modify the value of N for each run, based on the eleven point average precision for that run. In TREC-2, their experiments concentrated on methods of combining runs based on the similarity values of a document to each query for each of the runs. After some tests, the best choice is to weight each of the separate runs equally and not favor any individual run or method, but sometimes some runs has to be weighted more or less, depending on their performance. This method of merging the runs help the retrieval system to make a trade-off between the runs' errors. During the tests have been considered six different way to combine the runs:

- **CombMIN**: it is used to minimize the probability that a non-relevant document would be highly ranked;
- **CombMAX**: it is used to minimize the number of relevant documents being poorly ranked;
- **CombMED**: it is used to take the median similarity value (to solve the previous methods' problem) instead of taking a value from only one run;
- **CombSUM**: it is used to take the sum of the set of the similarity values;
- **CombANZ**: it is used to take the average of the non-zero similarity values, so it ignores the runs that fail to retrieve a relevant document;
- **CombMNZ**: it is used to consider the higher weights to documents retrieved by multiple retrieval methods.

We have to point out that the first two methods have a specific objective but they do not care about the possible problems that they can generate on the other retrieved documents. During the tests *CombMIN* has worse performance than all the single runs, on contrary *CombANZ* and *CombMNZ* methods have better performance than the individual runs, it is possible maybe because they produce the same ranked sequence for all the documents retrieved by all five individual runs.

3. Initial Attempts

Before going into the details of our final solution, it is useful to describe previous approaches we took into account to solve the problem and why we chose not to explore them further.

3.1. Parsing Documents

Multiple parsers have been developed to parse the documents from the provided collection. The most trivial parser, called *P1*, extracts the *sourceText* and *discussionTitle* elements from the corpus documents. The second parser *P2* extracts the elements related to the conclusion, the premise, the discussion title, and all the text from the *sourceText* field in between the premise and the conclusion. The third parser *P0*, which at the end we decided to use for our more advanced experiments, extracts only the discussion title, the premise and the conclusion for each document. Table 1 shows how the index statistics are affected by each of these parsers.

Table 1

Statistics regarding three indexes generated using the three implemented document parsers. The analyzer is always the same. Time Ratio is obtained by dividing the time taken by each parser by the time taken by the fastest parser.

Parser	Term Count	Storage (MB)	Time (seconds)	Time ratio
<i>P0</i>	1078017	195	99	1,00
<i>P1</i>	1196289	1745	507	5,12
<i>P2</i>	1153517	706	249	2,52

3.2. Query Expansion

When we have developed the software to create the index, we have deeply thought about how we could have used the resulting tokens from the analysis of the topic query.

3.2.1. OpenAI GPT-2

In the attempt of expanding the queries we came across the OpenAI GPT-2 model, a Machine Learning algorithm that generates synthetic text samples from an arbitrary input.

The idea was to use this powerful algorithm to make a query expansion, giving as input the topic title to generate a more complete phrase with hopefully new words that could help the searching part. Unfortunately the output of GPT-2 is not always what we expect. For example if we give it as input the tokenized query title, that could be only made of 2 words, the output is a not very useful dialog for our task. Another problem is the structure of the query, in fact

since they are all questions, the GPT-2 algorithm generates an answer for them which still is not what we were interested in. The problem persists also if we remove the question mark at the end of the phrase. In fact, the queries still have a question structure. For these reasons, we have decided to set aside this kind of approach.

3.2.2. Randomly Weighted Synonyms

An approach initially devised for query expansion, but which we later decided not to explore further, was to generate multiple queries for the same topic, each with randomly generated synonym boost values. For each query, the rankings of 1000 documents were then generated, and finally all the rankings were merged into one. The first performances obtained by this method did not encourage us to proceed with the development because there were many possible paths to follow from that point and the search time increased considerably.

3.3. Minimum body length

During the process of documents exploration, useful to detect which field are needed and present in each collection, we notice that in some documents the *sourceText* field were constituted by useless text without a single relevant information about its topic. In order to avoid this kind of documents ending up in the inverted file, we have tried to include, during the indexing process, a check on the length of the *ParsedDocument*'s body. If this field, the one that we have considered as the union of the *conclusion* and *premise* fields, is made up of less than a certain number of token, recurrent aspect that the parsing phase highlights for such instance, we avoid to consider them in the indexing phase.

After doing some tests with different values for the "min body length" (5, 10, 15), we have compared the results of this kind of solution with the results obtained without using it and we have discovered that, taking as example "min body length" equal to 10, the number of retrieved document switches from 48781 to 48764 and the number of relevant document switches from 1263 to 1257 (the other evaluation measures are no so affected by this change).

Considering this result we have decided not to use this kind of document pre-processing to avoid the discarding of some document that could be considered relevant in the qrels file (and so for an user) even if they don't seem like it.

3.4. Re-ranking with discussion ID

One of our primary goals is to improve the final ranking of the documents. With this assumption we tried to improve performance by using re-ranking.

As a first analysis we observed that, in the documents of the dataset, the posts related to the same discussion had the same first part of the document ID. Based on the assumption that only posts from some discussion are relevant to a query, we tried to index those posts as a single document. We finally obtained an index with a discussion-based clustering of documents. In the searching phase, we firstly searched the query in the normal index, retrieving the classic ranking of single documents. Secondly we searched the same query in the second index of discussion clusters, thus obtaining a ranking of discussions. Finally the scores of documents in the first ranking were increased based on the rank of their respective discussion in the second

ranking. Unfortunately this approach did not provide the desired results because it assumes that all posts related to a discussion have the same relevance to the searched topic. In fact, it was found that some posts do not contain any useful information to argue the searched topic but are part of a discussion that is really relevant.

3.5. OpenNLP attempt

Exploring new solutions, we have even tried to implement a version of the program that uses the OpenNLP Machine Learning toolkit in order to see which advantages a tokenization able to distinguish location, personal nouns and so on could provide for the solution.

Following this path we have encountered an error which requires significant changes in the workflow of what we had done up to that point. For this reason we have decided not to keep going on with this branch.

4. Methodology

The goal of this task is to retrieve relevant arguments from online debate portals, given a query on a controversial topic. We have checked the dataset and we have noticed that it is composed by five JSON files. We have also read some documents and we have noticed that the main structure is the same for each one but with some different fields. To use the documents with *Lucene* we had to parse the documents. To do so, we have used the *Jackson library* and we have implemented our parser *P0* that takes the premises, conclusions, the document title and the stance attribute (pro or cons) of the documents.

4.1. Indexing

We have built four different parts starting from the *Lucene* default ones: *ArgsParser*, *Shanks-Analyzer*, *DirectoryIndexer* and the *Searcher*. There is indeed a *ShanksTouche* class which has the main method and allows us to setup parameters for indexing and analyzing parts. After converting the documents into something that *Lucene* can work on, we focused ourselves on the development of how the indexer is created, in particular on how the analyzer module works. Into the tokenization phase we have used the *StandardTokenizer*, the *LowerCaseFilter*, the *EnglishPossessiveFilter* and the stop-words *StopFilter*.

The arguments of the collection are stored in the index using four fields: ID, Title, Body, and Stance (pro or con).

4.1.1. Custom Stop-List

As you can see in Figure 4 we have compared the baseline with the default stop-list and with our custom stop-list; after that comparison we have decided to use a custom stop-list to better achieve the project goal. Our stop-list contains 1362 words that are derived from the merging of other stop-lists (e.g. smart and lucene) typically used. Our custom stop-list reduces memory usage by approximately 38% and indexing time by almost 20%.

4.2. Searching

In the searching phase of our program we focus our attention in finding strategies to improve the general quality of the results: experimenting with different approaches like query expansion based on *WordNet* synonyms or defining queries to score differently the fields of the documents. The approach we have chosen to perform involves the use of *BooleanQuery*. In this way it is possible to assign specific weights to every term of the query (boosts).

Finally, we decided to use and test both *BM25Similarity* and *LMDirichletSimilarity* similarities and then their combination with a *MultiSimilarity* to get the document scores.

4.3. Re-Ranking

During various experiments we have noticed that some measures of similarity and set of parameters favored the precision at the expense of recall and vice versa. With the purpose of combining advantages of both cases, we opted for a re-ranking method which exploits different similarities and query parameters. Our implementation is made of two steps. In the first one we use a query able to obtain a higher recall value when searching the index for relevant documents, whose parameters and similarity are decided based on empirical trials. In the following step we use a second query with better performance in terms of *ndcg* to re-evaluate the returned documents and so re-ranking them according to the new score. We call *maxRecall* the first query and *maxNdcg* the second one. According to our implementation, this approach turns out to be effective on the 2020 topic set, but at the expense of the time spent in the search phase, which increases quite substantially.

5. Experimental Setup

Our work is based on the following experimental setups:

- Repository: <https://bitbucket.org/upd-dei-stud-prj/seupd2021-goldr>;
- During the develop and the experimentation we have used our own computer and in the end we have run our code using *Tira*;
- Evaluation measure: BM25, LMDirichlet and a "MULTI" where both were combined;
- *Apache Maven*, *Lucene*;
- Java JDK version 11;
- Version control system *git*.
- *trec_eval* tool [5]

The collection is a set of 387,740 arguments crawled from *debatewise.org*, *idebate.org*, *debatepedia.org*, *debate.org* and 48 arguments from Canadian parliament discussions. We used the 50 topics from Touché 2020 Task 1 [6] of the contest to train and refine our search engine. Furthermore we developed the source code collaborating through the BitBucket platform.

6. Results and Discussion

In this section we provide graphical and numerical results about the experiments we conducted during the development of the project. We also discuss these results to derive some useful insights.

6.1. Our Baseline

In order to be able to track performance progress during development, we created our own baseline. For each of the three parsers $P0$, $P1$, $P2$ we produced a run using a simple analyzer that uses a *StandardTokenizer*, *LowerCaseFilter* and a *StopFilter* where the stop-list used is the standard one offered by *Lucene*. The similarity used is BM25. Figure 1 compares the three baselines. From these results we decided to develop our approach based on the $P0$ parser. Figures 2 and 3 show how $P0$ compares to the other parser evaluating performance per topic. The choice of $P0$ was motivated not only by the statistics in Table 1 but also by its overall efficacy as shown in the following figures. Figure 1 shows how differently the three approaches can behave. $P0$, that extracts only the discussion title, the premise and the conclusion for each document (with its stance), highlights better performance with respect to the other two retrieving more relevant document across the entire run. Obviously we chose to discard $P1$ because his performances were inferior to the other two as shown in Figure 2. To better understand the behavior of $P2$ versus to $P0$, Figure 3 compare the *per topic average precision*. Approximately 10% of the topic, in particular topic 3, 23, 34 ad 43 to 46 tend to give nicer results with parser $P2$.

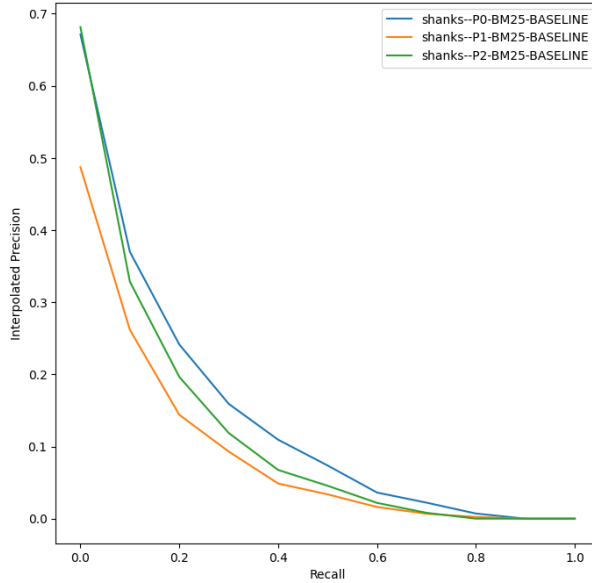


Figure 1: Plot of the Interpolated precision against recall of the three baselines.

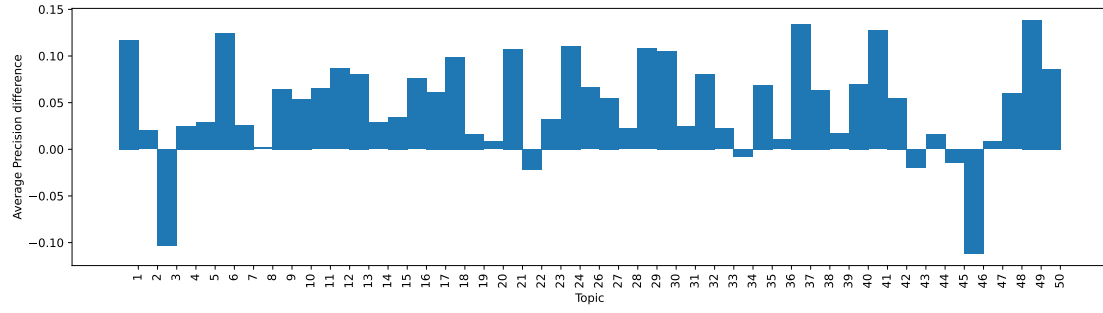


Figure 2: Per topic Average Precision Difference = $AP_{P0} - AP_{P1}$.

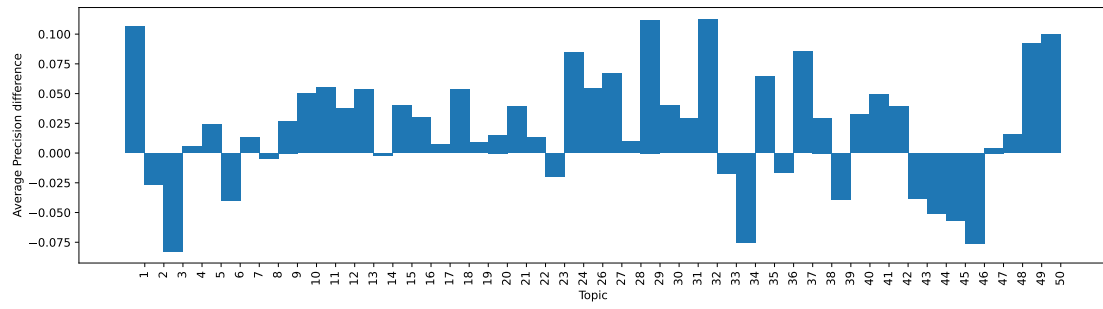


Figure 3: Per topic Average Precision Difference = $AP_{P0} - AP_{P2}$.

6.1.1. Baseline with the custom stop-list

The results obtained by our custom stop list, compared to the one offered by Lucene are comparable.

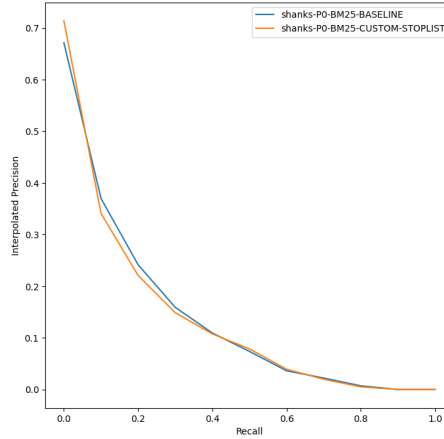


Figure 4: Plot of the Interpolated precision against recall of the two P0 baselines obtained by the Lucene stop-list and the custom stop-list.

6.2. Parameters Tuning

Our ultimate goal is to maximize the average value of $nDGC@5$ across all topics. To find the optimal parameter values, we performed extensive iterative tests, trying all combinations of parameters with values belonging to discrete intervals we defined. The same experiments were conducted using three different similarities : *BM25Similarity* (BM25), *LMDirichletSimilarity* (LMD), *MultiSimilarity* (MULTI). The MultiSimilarity we have used, combines BM25Similarity and LMDirichletSimilarity. The method we developed is governed by five parameters, when using re-ranking they become ten. Out of those ten, five parameters concern the query search in the index which aims to maximize the overall recall, the other five affect the re-ranking process and are chosen to maximize $nDGC@5$.

6.3. Optimal parameters

The optimal set of parameters we found for the *maxRecall* query and for the *maxnDCG* query are available in Tables 2 and 3. The best similarity measure to maximize $nDGC@5$, according to our empirical tests is LMD. To maximize recall, on the other hand, the best measure is MULTI. As it can be seen in Table 2, the best value of *tBoost* for *maxnDCG* is 0. This allowed us to come to the conclusion that considering the title of the discussion (which is the same for all posts in it) can be misleading with respect to the relevance of the content of each post that is part of it. The title, however, is very useful to obtain higher recall. This intuition is what pushed us to abandon the method of re-ranking based on discussion ID, which is described in 3.4.

Table 2

The optimal parameters obtained by training on the 2020 topics (topic description not considered).

Query	Similarity	tBoost	sBoost	pBoost	pDist
<i>maxRecall</i>	MULTI	3,50	0,15	1,75	12
<i>maxnDCG</i>	LMD	0,00	0,05	0,75	15

Table 3

The optimal parameters obtained by training on the 2020 topics (considering the topic description).

Query	Similarity	tBoost	sBoost	pBoost	pDist
<i>maxRecall</i>	MULTI	3,50	0,15	1,75	12
<i>maxnDCG</i>	LMD	0,30	0,05	1,75	15

6.4. maxnDCG and maxRecall

In this section we want to compare the two queries *maxnDCG* and *maxRecall* with the optimal parameter values established by the tests. From the graph in Figure 5 we can see that the precision is significantly higher for *maxnDCG*, while *maxRecall* has a better recall on the whole ranking. From these data we believe to have obtained the desired result from the two queries. In Figures 6, 7, 8 we compare the Average Precision per topic for the three test runs. We can notice how there is not much difference between maxRecall and P0. The same is not true for *maxnDCG*, the latter proves to be better than P0 in almost all topics. This scenario is further confirmed by Figure 8, which shows the dominance of $AP_{maxnDCG}$ over $AP_{maxRecall}$. The performance can be further compared with the numerical results reported in Table 4. From these data we can realize the actual trade-off between the two approaches. The advantage in terms of recall for *maxRecall* is significant compared to *maxnDCG*. The same advantage applies in the opposite way for precision and *nDCG*.

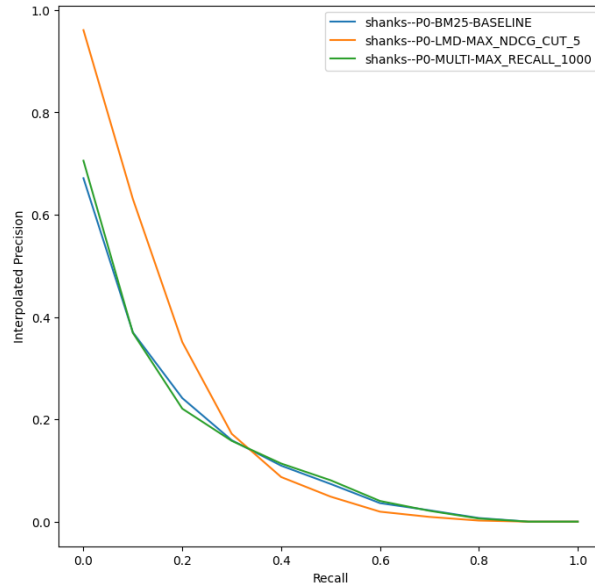


Figure 5: Plot of the Interpolated precision against recall for the *maxnDCG* query versus *maxRecall*.

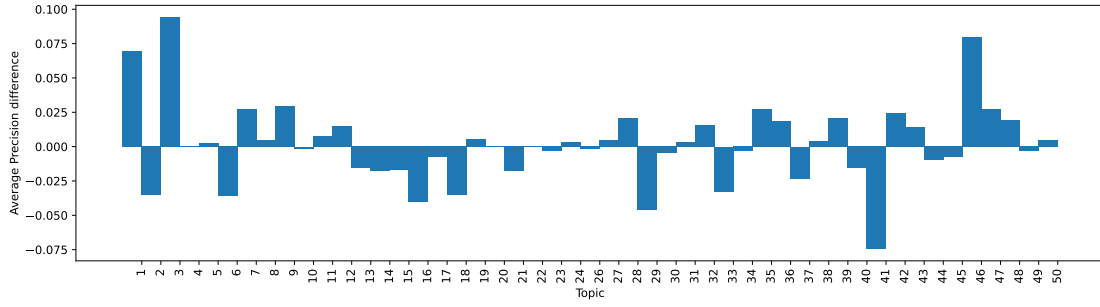


Figure 6: Per topic Average Precision Difference = $AP_{maxRecall} - AP_{p0}$.

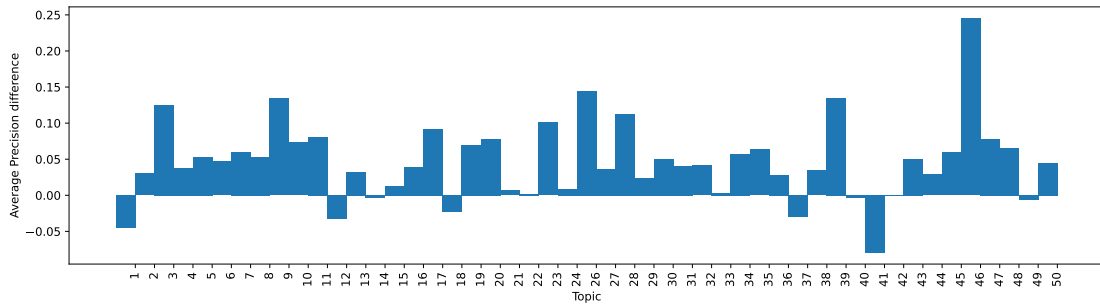


Figure 7: Per topic Average Precision Difference = $AP_{maxnDCG} - AP_{p0}$.

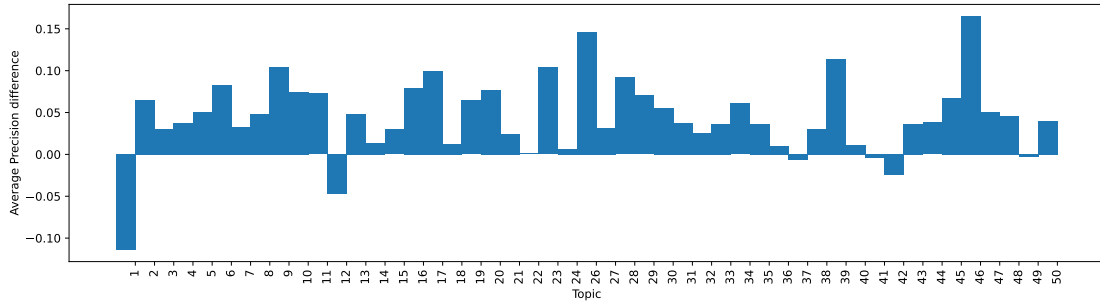


Figure 8: Per topic Average Precision Difference = $AP_{maxnDCG} - AP_{maxRecall}$.

6.5. Re-Ranking Results

Here we compare the *Re-Ranking* approach described in 4.3, with the previous ones. Figure 9 shows how re-ranking based on *maxNdcg* (in red), greatly improves the interpolated precision of *maxRecall* (in green). It is also possible to note that it allows for better performance with

respect to $maxNdcg$ alone (in orange). Figure 10 shows the Average Precision value obtained by re-ranking for each topic. From it, it is possible to identify the most problematic topics, for which the method developed is not very effective. In particular, the most critical topics are: {2, 8, 22, 40, 44}. From Figure 11 we can see that the *Re-Ranking* method improves on almost every topic w.r.t the baseline. In Figure 12 it can be seen that the *Re-Ranking* method improves on $maxNdcg$ on many topics, except for topic 10 and 20. Figure 13, as predictable, clearly shows the better performance achieved by *Re-Ranking* w.r.t. $maxRecall$, only the first topic is penalized by the re-ranking process.

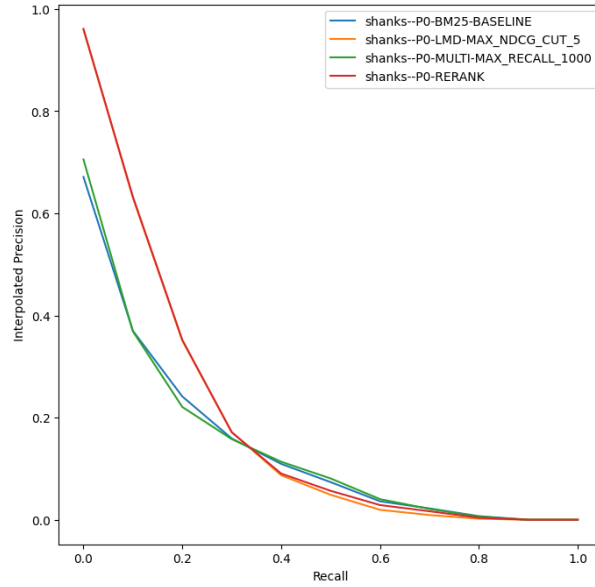


Figure 9: Plot of the Interpolated precision against recall comparing baseline, $maxRecall$, $maxnDCG$, *Re-Ranking*

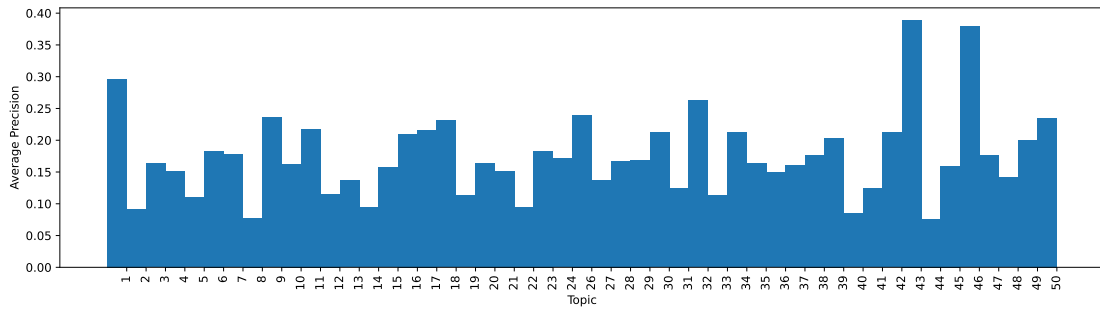


Figure 10: Per topic Average Precision for the Re-Ranking approach.

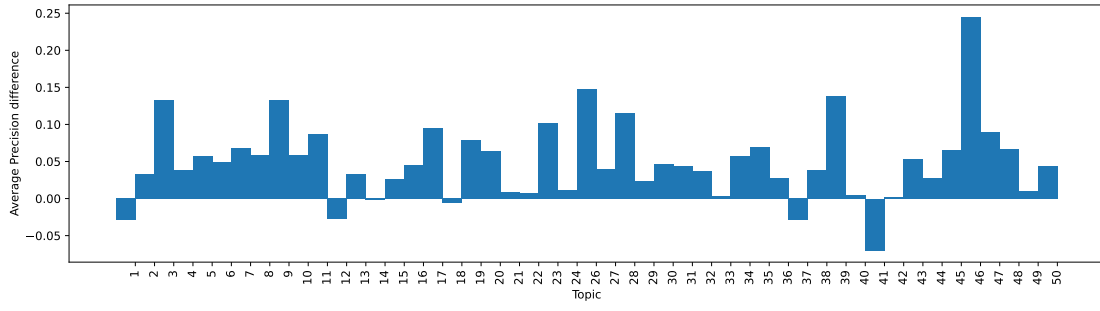


Figure 11: Per topic Average Precision Difference = $AP_{re-ranking} - AP_{p0}$.

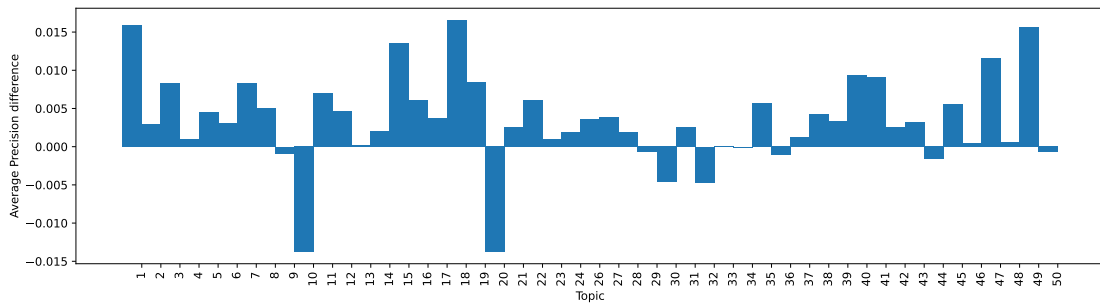


Figure 12: Per topic Average Precision Difference = $AP_{re-ranking} - AP_{maxnDCG}$.

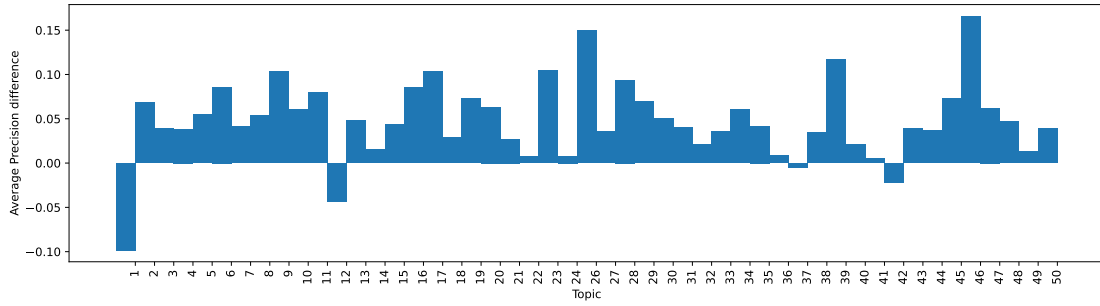


Figure 13: Per topic Average Precision Difference = $AP_{re-ranking} - AP_{maxRecall}$.

Table 4

Some numerical results for comparing performance of *Re-Ranking*, *maxnDCG*, *maxRecall* and the baseline.

RUN	num_rel_ret	map	P_5	recall_1000	nDCG	nDCG@5
P0-RERANK	1263	0.1750	0.6490	0.6631	0.4979	0.5495
P0-LMD-MAX_nDCG@5	1103	0.1717	0.6490	0.5803	0.4764	0.5495
P0-MULTI-MAX_RECALL_1000	1263	0.1276	0.4082	0.6631	0.4283	0.3117
P0-BM25-BASELINE	1250	0.1256	0.3796	0.6563	0.4216	0.2871

6.6. RUN Submission

The five run we decided to submit are the following:

- run-1 : *Re-Ranking* approach.
- run-2 : like run-1 but proximity searches are only with pairs of subsequent tokens.
- run-3 : *maxnDCG* query with LMDirichletSimilarity
- run-4 : *maxnDCG* query with MultiSimilarity
- run-5 : *maxRecall* query with MultiSimilarity

Table 5 shows the numerical results we obtained for each run. Out of all, the first run is the best one overall.

Table 5

Numerical statistics from the trec-evaluations of the 5 run on the 2020 topic set.

RUN	num_rel_ret	map	P_5	recall_1000	nDCG	nDCG@5
shanks-run-1	1263	0.1750	0.6490	0.6631	0.4979	0.5495
shanks-run-2	1255	0.1735	0.6408	0.6588	0.4960	0.5432
shanks-run-3	1103	0.1717	0.6490	0.5803	0.4764	0.5495
shanks-run-4	1199	0.1497	0.4816	0.6312	0.4536	0.3866
shanks-run-5	1263	0.1276	0.4082	0.6631	0.4283	0.3117

SUBMISSION UPDATE:

NEW EXPERIMENTAL RESULTS AFTER A NEW CORRECTED VERSION OF THE QRELS WAS RELEASED

All the previous results are based on an incorrect version of the .qrels file for the 2020 topics. Since we only knew about the new corrected version when the deadline was approaching, we could not recreate all the graphs and comparisons in 6. However, we managed to find the new optimal parameter sets and repeat the five run.

The new data is provided in Tables 6 and 7.

Table 6

The optimal parameters obtained by training on the 2020 topics (WITH CORRECTED QRELS).

Query	Similarity	tBoost	sBoost	pBoost	pDist
maxRecall	<i>MULTI</i>	0,3	0,2	0,75	12
maxnDCG	<i>LMD</i>	0,15	0,05	0,75	17

Table 7

Numerical statistics from the 5 run on the 2020 topic set (WITH CORRECTED QRELS).

RUN	num_rel_ret	map	P_5	recall_1000	nDCG	nDCG@5
shanks-run-1	795	0.3146	0.6245	0.8705	0.6521	0.6407
shanks-run-2	790	0.3126	0.5959	0.8671	0.6521	0.6213
shanks-run-3	770	0.3141	0.6245	0.8502	0.6479	0.6407
shanks-run-4	788	0.2546	0.4327	0.865	0.5839	0.4391
shanks-run-5	795	0.2565	0.4449	0.8705	0.588	0.4513

7. Statistical Analysis

Here we analyze our models with some important statistics to evaluate them in a deeper way via hypothesis testing. We used *ANOVA*, *tStudent* test and through *boxplots*. All the analyses focused on the mean of two key metrics, average precision and nDCG@5. We analyzed all 5 different retrieval models, the ones described in the section 6.6.

A boxplot gives a visual representation about location, symmetry of the data, dispersion and presence of outliers (points that escape the construction of the boxplot).

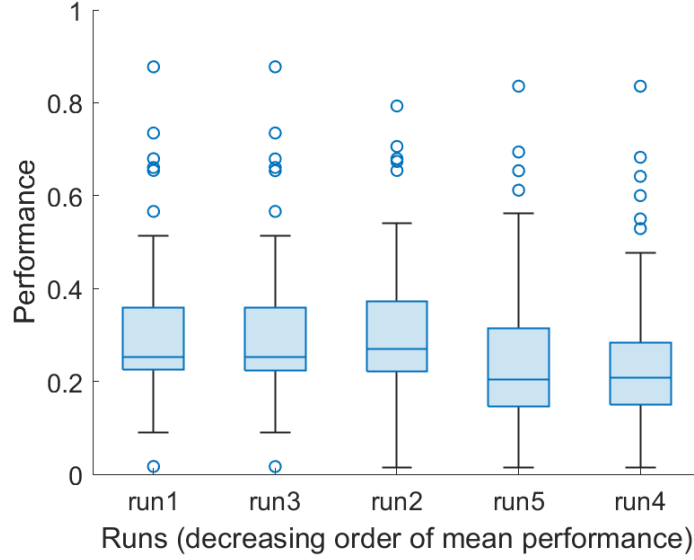


Figure 14: Boxplots of the 5 runs of Aveage Precision

It can be appreciate that *run1*, *run3* have almost an identical structure: the median, interquartile range (IRQ) and whiskers length. On the other hand, *run4* and *run5* exhibit less performance in this metric because those runs were tuned for maximize recall or nDCG@5, the bottom whiskers in fact are closer to zero meaning that for some topics the system did not retrieve enough relevant documents. All the runs have outliers, the points above the whiskers, represented by circles, are topics which perform better than the others, or worse if they are below the boxplot. *Run1* and *run3* are skewed to the right and show less variance compared to other systems.

Looking at the boxlot of nDCG@5 reveals the same behavior seen before, *run1* and *run3* produce higher scores compared to other runs and they seem identical in performances. *Run3* has better results w.r.t *run4* so we can conclude that using different similarities change dramatically the results. *Run2* shows less IRQ among others, in particular compared to *run1* that share the same architecture with the only difference in the proximity parameter. We can say that *run1* is able to score higher score exploiting the more flexible proximity parameter.

Further analysis with anova and tTest will help to understand the possible similarities or not between the systems.

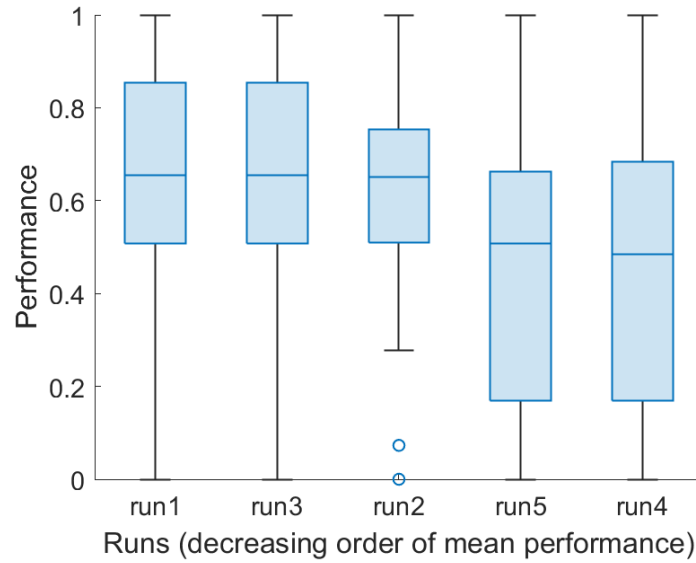


Figure 15: Boxplots of the 5 runs of nDCG@5

7.1. Hypothesis testing

The first tool that we utilize is *ANOVA (Analysis of Variance)* a statistical test of whether or not the means of several groups are equal. H_0 , the null hypothesis that all the means are equal, is tested against the possibility to reject or don't reject it. As we can see in Figure 16 there are multiple factors to be taken into account to perform a correct analysis of the F-statistic. The system sum of squares, SS_{system} and the SS_{error} are divided by their degree of freedom to obtain the mean squares (MS). The F-statistic (F) is equal to the ratio of MS_{System} and MS_{Error} . Having a $p-value = 0.1849$ we cannot reject H_0 . To do that we should have had a value lower than α . The significance level α is set at 0.05.

We could conclude that our systems are statistically similar in mean average precision but performing the *ANOVA2* test the situation is reversed.

	1	2	3	4	5	6
1	'Source'	'SS'	'df'	'MS'	'F'	'Prob>F'
2	'Columns'	0.1994	4	0.0499	1.5630	0.1849
3	'Error'	7.6556	240	0.0319	[]	[]
4	'Total'	7.8550	244	[]	[]	[]

Figure 16: Anova1 results

	1	2	3	4	5	6
1	'Source'	'SS'	'df'	'MS'	'F'	'Prob>F'
2	'Columns'	0.1994	4	0.0499	34.8793	7.0038e-22
3	'Rows'	7.3812	48	0.1538	107.5816	1.6047e-115
4	'Error'	0.2744	192	0.0014	[]	[]
5	'Total'	7.8550	244	[]	[]	[]

Figure 17: Anova2 results

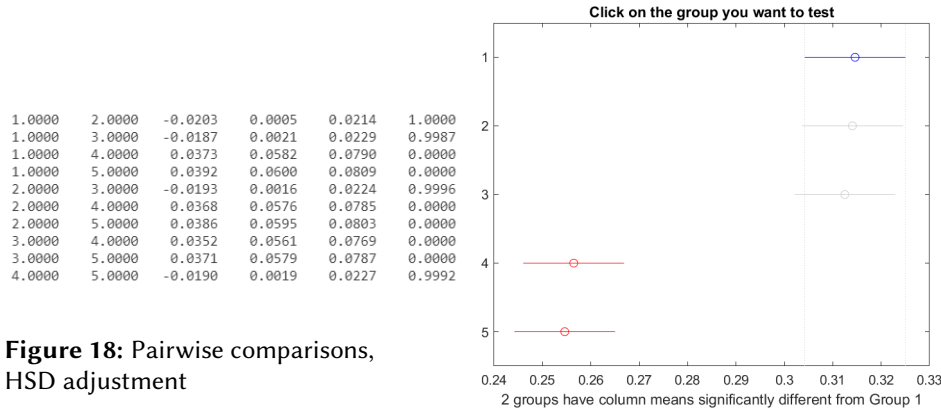


Figure 18: Pairwise comparisons, HSD adjustment

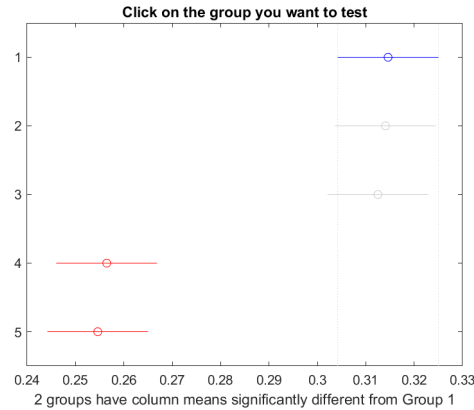


Figure 19: Multi comparison chart for AP

As it can be seen from Figure 18 and Figure 19, the null hypothesis can be rejected because the runs have statistically significant differences. The topic effect, that can be read in cell [3,2] of Figure 17, is able to express much more variance, that means that this variability is greater than the one of the systems as we can expect by topics.

Run 3 vs run 4: as we pointed out previously, the use of LMDirichlet similarity improve the results. For this pair the null hypothesis can be rejected. As we will see for the nDCG@5 analysis *run 4* and *run 5* are different from the others, and the *p-value* of their statistic suggest us that the runs are equivalent in mean. Instead, for the others (*run 1*, *run 3*, *run 2*) we cannot reject H_0 .

Unfortunately we have not been able to simplify the reading of the images. Here is a sort of conversion between the numbers on the axes and the real ones.

1 - 2 - 3 - 4 - 5 (image numbers) \rightarrow run1 - run3 - run2 - run5 - run4 (true values). This is applicable to Figure 18, 19, 21, 23

Moving the study to nDCG@5, in ANOVA table, Figure 20, the *p-value* is lower than the confidence level so we can reject the null hypothesis and claim that there is at least one mean between the various runs that differs significantly from the others. This table doesn't tell us which systems are different on average, it just tells us that there is at least one.

Tukey Honestly Significant Difference (HSD) test, as we did in AP case, answer to that point creating confidence intervals for all pairwise differences between the systems we want to compare, while controlling the family error rate. Otherwise, the probability of a type I(one) error would be magnified.

Run4 and *run5*, Figure 22, are statistically different from the others, while for them we cannot reject the null hypothesis. It seems that supporting the *MaxRecall* or *MaxnDCG* query approach does not yield performance benefits except through different similarity as in case of *run3* or re-rank(*run1* and *run2*). Accordingly to what the boxplot chart suggested for *run3* and *run4* we can reject H_0 , the p-value is below $\alpha = 0.05$. The different similarity approach is visible.

We can conclude that our *Re-ranking* approach is much more significant than a standard technique even if *run3* through this analysis it returns comparable results. In the future this behavior can be investigated carefully.

	1	2	3	4	5	6
1	'Source'	'SS'	'df'	'MS'	'F'	'Prob>F'
2	'Columns'	2.1162	4	0.5290	7.6877	7.6237e-06
3	'Error'	16.5160	240	0.0688	[]	[]
4	'Total'	18.6321	244	[]	[]	[]

Figure 20: nDCG@5 - Anova1 results

1.0000	2.0000	-0.1446	0	0.1446	1.0000
1.0000	3.0000	-0.1252	0.0194	0.1640	0.9962
1.0000	4.0000	0.0448	0.1894	0.3339	0.0033
1.0000	5.0000	0.0570	0.2016	0.3461	0.0013
2.0000	3.0000	-0.1252	0.0194	0.1640	0.9962
2.0000	4.0000	0.0448	0.1894	0.3339	0.0033
2.0000	5.0000	0.0570	0.2016	0.3461	0.0013
3.0000	4.0000	0.0254	0.1700	0.3145	0.0117
3.0000	5.0000	0.0376	0.1822	0.3267	0.0053
4.0000	5.0000	-0.1324	0.0122	0.1568	0.9994

Figure 21: nDCG@5 - Multiple pairwise comparison of the group means, *first column are the pairs of systems, the last one is the related p-value*

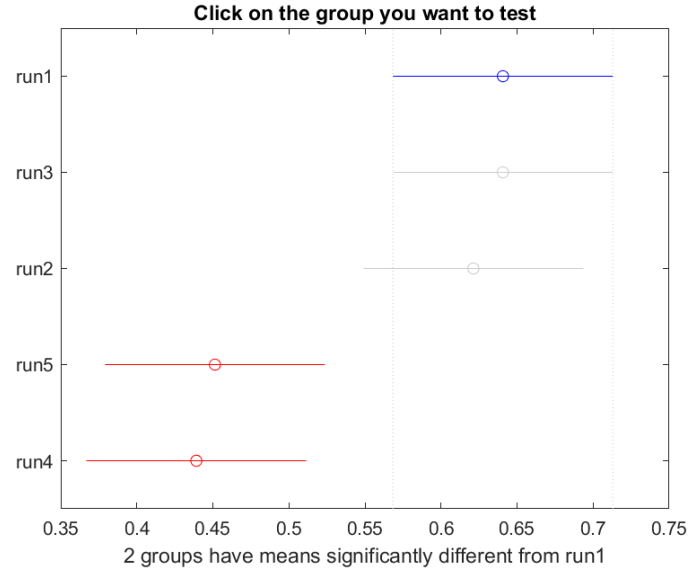


Figure 22: Multiple comparison chart for nDCG@5

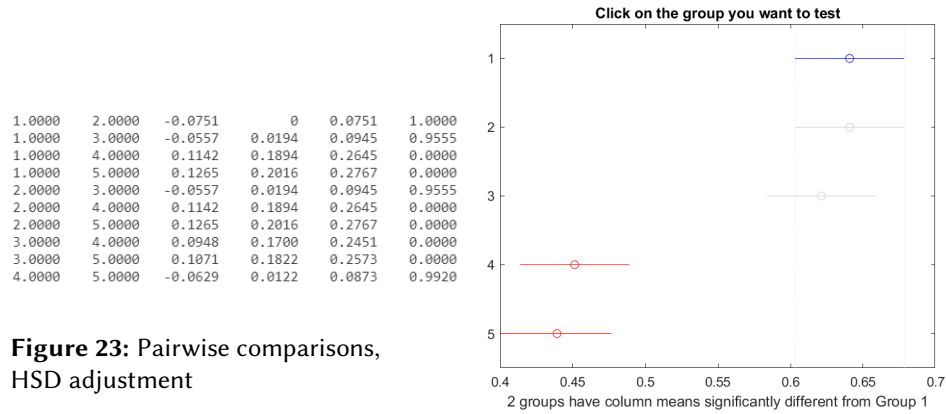


Figure 23: Pairwise comparisons, HSD adjustment

Figure 24: Multi comparison chart for nDCG@5

8. Failure analysis

Looking at the results obtained from the `trec_eval` program execution with our runs as a parameter we can see the performance of our information retrieval systems. In particular we are now focusing on finding and understanding for which topic the systems fail in achieving good performance. Therefore we have decided to apply this kind of evaluation to our best run (*shanks-run 1*). Looking at it, and using the *map* field as reference parameter for the following analysis, we have discovered that the top (*map*) performance come from the topics 42, 43 and 1 and the worst from 22, 12 and 44. Searching for the reason why, we understand that the main

weakness of the process is the lacking of an argument quality evaluation phase and of a more consistent lexical analysis process.

As an example we take and compare the topics 1, 12 and 44.

Topic 1: Should teachers get tenure?

Narrative: highly relevant arguments make a clear statement about tenure for teachers in schools or universities. Relevant arguments consider tenure more generally, not specifically for teachers, or, instead of talking about tenure, consider the situation of teachers' financial independence.

The process of stopwording applied to the topics lead to obtain the parsed query "teachers tenure", that even without the whole phrase construction explain very well what we are searching for. The document about this topic are well retrieved by our system, in fact if we look for example at the first 5 position (the most relevant ones for a browser) we can see from the comparison with the qrels that all of them are considered as "highly relevant".

Topic 12: Should birth control pills be available over the counter?

Narrative: highly relevant arguments argue for or against the availability of birth control pills without prescription. Relevant arguments argue with regard to birth control pills and their side effects only. Arguments only arguing for or against birth control are irrelevant.

The process of stopwording applied to the topics lead to obtain the parsed query "birth control pills available". This seems a quite explicative phrase but in the retrieval phase the system has trouble in finding the proper relevant results. The critical issue with this query is the distinction between high and low quality argument. The lacking of an effective argument quality process leads the program in failing the document quality evaluation. Due to this fact the system is unable to put in the right ranking position the appropriate document. We can notice this aspect even looking at the other topic fields (not only the map), in fact if we look for example at the growing of the recall parameter we can notice that is mainly focused in the tail of the process (when many documents are retrieved).

Topic 44: Should election day be a national holiday?

Narrative: highly relevant arguments explain why making election day a holiday is a good idea, or why not. Relevant arguments mention the fact or its remedy as one of the problems that elections have.

The process of stopwording applied to the topics lead to obtain the parsed query "election national holiday". The results obtained from this kind of search are quite bad in fact the system retrieves as relevant some discussion like "Potato day should be a national holiday" or "Star Wars day should be a national holiday". These mismatches came from the fact that the system does not recognize "election day" as unique mandatory query term and so document about similar topic that differ only for some word are wrongly retrieved.

To overcome these kind of issues it could be useful to equip the system with an argument quality evaluator and a sort of word pattern recognizer that catches the words which must not be separate and a way to identify structure like "subject, predicate, object" that assign higher weights to the subject and marks it as mandatory word in the documents to be retrieved.

9. Conclusions and Future Work

At the end of this experiment we have discovered that the changes we made, lead to obtain fairly good improvement in respect to the starting baseline [Table 4] of our information retrieval system. All the statistics have undergone a significant increase, as an example we can see the improvement as follows: MAP +39.3%, P5 +68.8%, nDCG@5 +91.4%.

Comparing our results with the last year ones we have noticed that they follows their trend for the nDCG@5 parameter, but looking at the applied strategies they seem to be fairly different. Our results are in line with those of last year.

The whole process we followed highlighted a lot of possible expansions that with more time could be implemented and explored. As an example it could be interesting the application of some sort of location word detection, personal nouns recognition and compound words discernment and an improvement or even the addition of a new ranking phase that allows the insertion of some kind of argument quality analysis. The possibility of implementing a different approach based on some kind of machine learning model remains open. As we stated in the initial attempt section 3 we dropped the idea of using GPT-2 because it returned us unsatisfying results, so in the future we can go more into details of this modern tool and improve in this way the performance of our solution.

Another possible approach that can be exploited derives from the paper [4]. It can be interesting because it takes different runs in input and combines them to reach the best result.

References

- [1] C. Fellbaum, WordNet: An Electronic Lexical Database, Bradford Books, 1998.
- [2] Y. Ajjour, H. Wachsmuth, J. Kiesel, M. Potthast, M. Hagen, B. Stein, args.me corpus, 2020. URL: <https://doi.org/10.5281/zenodo.3734893>. doi:10.5281/zenodo.3734893.
- [3] A. Bondarenko, M. Fröbe, M. Beloucif, L. Gienapp, Y. Ajjour, A. Panchenko, C. Biemann, B. Stein, H. Wachsmuth, M. Potthast, M. Hagen, Overview of Touché 2020: Argument Retrieval, CEUR-WS 2696 (2020).
- [4] E. A. Fox, J. Shaw, Combination of Multiple Searches, in: D. K. Harman (Ed.), The Second Text REtrieval Conference (TREC-2), National Institute of Standards and Technology (NIST), Special Publication 500-215, Washington, USA, 1993, pp. 243–252.
- [5] C. L. A. Clarke, N. Craswell, I. Soboroff, Overview of the TREC 2009 Web Track, in: E. M. Voorhees, L. P. Buckland (Eds.), The Eighteenth Text REtrieval Conference Proceedings (TREC 2009), National Institute of Standards and Technology (NIST), Special Publication 500-278, Washington, USA, 2010.
- [6] A. Bondarenko, M. Fröbe, M. Beloucif, L. Gienapp, Y. Ajjour, A. Panchenko, C. Biemann, B. Stein, H. Wachsmuth, M. Potthast, M. Hagen, Overview of Touché 2020: Argument Retrieval, in: L. Cappellato, C. Eickhoff, N. Ferro, A. Névél (Eds.), Working Notes Papers of the CLEF 2020 Evaluation Labs, volume 2696 of *CEUR Workshop Proceedings*, 2020. URL: <http://ceur-ws.org/Vol-2696/>.