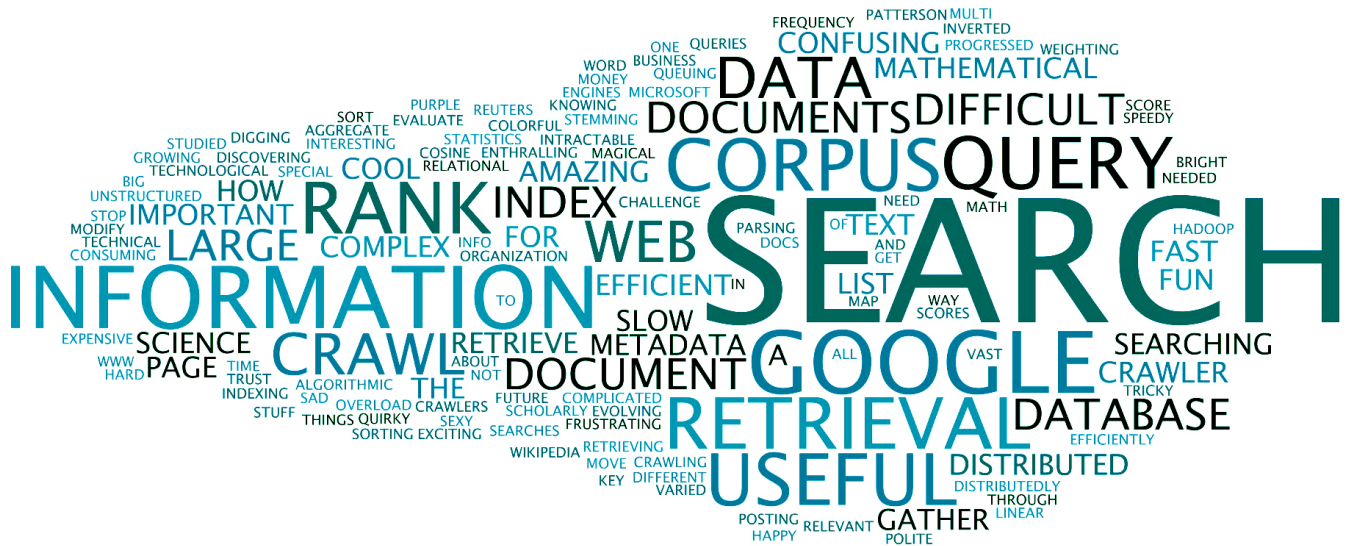


Information Retrieval and Data Mining

Porject 1 2020-2021



Team 2

Karypidis Stathis [Github](#)

Anagnostou Pantazis [Github](#)

Abstract: In this project we conduct an experiment on the basics of information retrieval. The project was based on the Lemur Project or Indri Search Engine. More specifically, we have 4 collections of texts (fbis, fr94, ft, latimes) and 150 information needs/topics which are tested on two different retrieval models. The first one is the default model that Indri has and its results will be our baseline results. On the other model we use query expansion with Relevance Feedback of terms of the collections or terms from a thesaurus. The version of the Lemur Project used is 5.11 on Ubuntu 16.04 and the evaluation tool used is trec_eval 9.07.

The code of this project can be found on github: <https://github.com/pantanag/IR-DM-Exercise-2020-21-Team-2>

Part A

In the first part of this project we use the default retrieval model of Indri. After the setup of Indri is complete, we use the command to build our index of terms of our collections (default tokenizer and Krovetz stemmer) after we have modified the parameter file "IndriBuildIndex.parameter.file.Example" to our needs. More specifically, we add the paths to our collections and the path that our index is going to be stored.

```
1 <parameters>
2 <corpus>
3   <path>/home/stathis/Downloads/IR-2019-2020-Project-1/ft</path>
4   <class>trectext</class>
5 </corpus>
6 <corpus>
7   <path>/home/stathis/Downloads/IR-2019-2020-Project-1/fr94</path>
8   <class>trectext</class>
9 </corpus>
10 <corpus>
11   <path>/home/stathis/Downloads/IR-2019-2020-Project-1/fbis</path>
12   <class>trectext</class>
13 </corpus>
14 <corpus>
15   <path>/home/stathis/Downloads/IR-2019-2020-Project-1/latimes</path>
16   <class>trectext</class>
17 </corpus>
18 <index>/home/stathis/Downloads/IR-2019-2020-Project-1/indices</index>
19 <memory>4096M</memory>
20 <storeDocs>true</storeDocs>
21 <stemmer><name>Krovetz</name></stemmer>
22 </parameters>
```

Next we use the command mentioned before `IndriBuildIndex IndriBuildIndex.parameter.file.Example` to build our index. With the command `dumpindex /path/to/index v` we can see our index too. The next step is the creation of a joint file containing all the qrels files (relevant judgments) with the use of the command `cat qrels.301-350.trec6.adhoc qrels.351-400.trec7.adhoc qrels.401-450.trec8.adhoc > qrels.301-450.trec.adhoc` and the creation of a file containing all the information needs (topics) with the corresponding command `cat topics.301-350.trec6 topics.351-400.trec7 topics.401-450.trec8 > topics.trec`. If we are given the above files we can skip to next part which is the creation of the three (3) queries files. After having read the queries file with only the titles (was given as an example) we proceed to remove any punctuation and handle any specific characters so that our 3 files are in the best format possible (and that Indri can understand them).

Our results were the below:

Open	Save	Open	Save	Open	Save
eval_results_titles.txt (~Downloads/trec_eval-9.0.7) - gedit		eval_results_titles-desc.txt (~Downloads/trec_eval-9.0.7) - gedit		eval_results_titles-desc-narr.txt (~Downloads/trec_eval-9.0.7) - gedit	
runid	all	runid	all	runid	all
num_q	all	num_q	all	num_q	all
num_ret	all	num_ret	all	num_ret	all
num_rel	all	num_rel	all	num_rel	all
num_rel_ret	all	num_rel_ret	all	num_rel_ret	all
map	all	map	all	map	all
gm_map	all	gm_map	all	gm_map	all
rprec	all	rprec	all	rprec	all
bpref	all	bpref	all	bpref	all
recip_rank	all	recip_rank	all	recip_rank	all
tprec_at_recall_0.00	all	tprec_at_recall_0.00	all	tprec_at_recall_0.00	all
tprec_at_recall_0.10	all	tprec_at_recall_0.10	all	tprec_at_recall_0.10	all
tprec_at_recall_0.20	all	tprec_at_recall_0.20	all	tprec_at_recall_0.20	all
tprec_at_recall_0.30	all	tprec_at_recall_0.30	all	tprec_at_recall_0.30	all
tprec_at_recall_0.40	all	tprec_at_recall_0.40	all	tprec_at_recall_0.40	all
tprec_at_recall_0.50	all	tprec_at_recall_0.50	all	tprec_at_recall_0.50	all
tprec_at_recall_0.60	all	tprec_at_recall_0.60	all	tprec_at_recall_0.60	all
tprec_at_recall_0.70	all	tprec_at_recall_0.70	all	tprec_at_recall_0.70	all
tprec_at_recall_0.80	all	tprec_at_recall_0.80	all	tprec_at_recall_0.80	all
tprec_at_recall_0.90	all	tprec_at_recall_0.90	all	tprec_at_recall_0.90	all
tprec_at_recall_1.00	all	tprec_at_recall_1.00	all	tprec_at_recall_1.00	all
P_5	all	P_5	all	P_5	all
P_10	all	P_10	all	P_10	all
P_15	all	P_15	all	P_15	all
P_20	all	P_20	all	P_20	all
P_30	all	P_30	all	P_30	all
P_100	all	P_100	all	P_100	all
P_200	all	P_200	all	P_200	all
P_500	all	P_500	all	P_500	all
P_1000	all	P_1000	all	P_1000	all

The metrics we are interested in are the below:

Metric	title	title+desc	title+desc+narr
Mean Average Precision	22.17%	22.95%	22.4%
R-Precision	26.91%	27.67%	27.01%
Precision@10	42.73%	43.67%	43.33%

While observing the results we can see that our metrics dont vary so much across the three fields. Obviously the queries containing both the title and the desc give the best results (based on the metrics above). While adding the desc clearly boosts our results, adding the narr field too doesnt improve our search. This indicates that adding the desc field provides us with extra information but adding the narr field too comes to surplus. Let's look at the arithmetic mean of the Average Precision or MAP (mean of AP at the different values of Recall or area below the curve Precision - Recall). As we can see its between 22-23% for the first 3 cases. Then we examine the R-precision or the percentage $\frac{r}{R}$ where **r** is the number of our relevant results of a query until the first **R** texts where **R** is the number of relevant texts for the query. This metric has also not so good results being between 26-27%. Finally, regarding the Precision@10 we notice that it greatly higher than our previous two metrics. The results of the MAP and R-precision indicate that our retrieval quality isn't so good while the results of the Precision@10 indicate that even for the first 10 texts, 4 or 5 max texts are relevant with our search. In the next part we examine the query expansion with Relevance Feedback.

Part B

As mentioned before, in this part of the project, we use query expansion. Moreover, we acquire the 20-most frequent terms on the 15 most relevant texts of our query (based on the results of the Relevance Feedback). Apart from that we use a thesaurus to produce synonyms of those terms and add them to our improved query too. To do this, we must first examine our results to get the 15 most relevant texts/titles (files: "**results-titles.trec**" and "**results-titles-desc.trec**"). The function responsible for that is the following:

```
1 def get_top15(filename):
2     #Open file with results
3     f = open(filename+".trec","r")
4     lines = f.readlines()
5     previous = ""
6     num_queries = []
7     final = {}
8     counter = 0
9     #Read every line and check for a new number query
10    for i in range(0,len(lines)):
11        line = lines[i]
12        line_values = line.split()
13        if(line_values[0] != previous):
14            num_queries.append(line_values[0])
15            previous = line_values[0]
16            #New query found so start keeping best 15
17            start = True
18            if(start):
19                temp = []
20                for k in range(0,15):
21                    temp.append(lines[i].split()[2])
22                    i = i+1
23                start = False
24                final[num_queries[counter]] = temp
25                counter = counter +1
26    return final
```

After we have found the 15 most relevant titles, we must find the corpus its title belongs to. To do this we built a lexicon (stored as a json file) so that we have a correspondence between the texts (titles) and the corpus they belong to. The file containing them is the **paths.json** and its necessary for the **partB.py** to run (must be in the same folder). In case we don't have the lexicon then the function searches all the collections to find the right one and then searches the right text (the use of the lexicon is highly recommended because it's a lot faster).

After we have found the right collection the article belongs to all that remains is for us to open it and extract the text. Once we have extracted the text, we tokenize it and we use the library **Counter** to keep track of the frequency of its term in the text.

```
1 #Function to retrieve text from the file obtained before
2 def retrieve_text(filename,title):
3     text = ""
4     #Pattern for html and etc tags
5     cleanTags = re.compile('<.*?>')
6     #Open and read file
7     with codecs.open(filename,"r",encoding="latin-1") as data:
8         lines = data.readlines()
9         for i in range(0,len(lines)):
```

```

10         line = lines[i]
11         #Find the right section of our file and gather only the <TEXT> part
12         if title in line:
13             i = i + 1
14             while("<TEXT>" not in lines[i]):
15                 i=i+1
16             i = i+1
17             while("</TEXT>" not in lines[i]):
18                 #Remove html tags and dont add empty or newlines
19                 tempLine = re.sub(cleanTags,"",lines[i])
20                 if tempLine not in ["","\n"]:
21                     text = text + tempLine
22                 i=i+1
23         #Remove punctuation
24         out = text.translate(text.maketrans('','',string.punctuation))
25         return out
26 #Create dictionary of tokenized words and their frequency
27 def make_dict(text):
28     ks = krovetz.PyKrovetzStemmer()
29     #Use nltk tokenizer to get every word in the text
30     lista = nltk.word_tokenize(text)
31     lista = [ks.stem(word) for word in lista]
32     counter = Counter(lista)
33     #Create stopwords array to dismiss them later
34     stopwords = nltk.corpus.stopwords.words('english')
35     stopwords.extend(list(string.ascii_lowercase))
36     stopwords.extend([word.capitalize() for word in stopwords])
37     pops = set(stopwords).intersection(counter.keys())
38     for i in pops:
39         counter.pop(i)
40     most_common = sorted(counter.items(),key = lambda pair: (-pair[1],pair[0]))
41     #Return only first 50 words
42     return most_common[0:50]

```

As we can see in the above code our function returns the 50 (not 20) most frequent terms in the 15 most relevant texts that were retrieved. The reason behind this is that we want to avoid cases of double-terms when we proceed to query expansion. Finally, all that's left to do is the query expansion while avoiding the double effect mentioned before and cases of vehicle-vehicles (by using **Krovetz Stemmer**). Our improved queries are then built and ready to be tested. The function responsible for all those things mentioned is the below:

```

1 def create_improved_queries(json_file,originalQueries,output,weight):
2     #Data extraction from json file
3     f = open(json_file+".json",'r')
4     data = json.load(f)
5     f.close()
6     keys = sorted(data.keys())
7     currPath = os.getcwd()
8     if not os.path.exists("Improved Queries"):
9         os.mkdir("Improved Queries")
10    os.chdir(currPath + "/Improved Queries")
11    #Create dictionary for current (old) queries
12    count = 301
13    queries = {}
14    originalQueries = remove_pun(originalQueries)
15    for q in originalQueries:
16        queries[str(count)] = q
17        count+=1

```

```

18 #Create improved queries
19 improvedQueries = []
20 for key in keys:
21     temp = ""
22     items = data[key]
23     #Count to keep track of number of iterations (we want top 20 only)
24     tempCount = 0
25     for item in items:
26         #Check for duplicates
27         if item[0].lower() not in queries[key].lower() and tempCount < 20:
28             temp = temp + " " + item[0]
29             tempCount +=1
30     #Build new improved query
31     improvedQueries.append(temp)
32 write_queries(originalQueries,output+"extra-20-queries",improvedQueries,weight)
33 queries_extra = combine_texts(originalQueries,improvedQueries)
34 tempCount = 0
35 #Create synonyms queries
36 synQueries =[]
37 ks = krovetz.PyKrovetzStemmer()
38 for query in queries_extra:
39     temp = ""
40     tokenized = nltk.word_tokenize(query)
41     words = [ks.stem(token).lower() for token in tokenized]
42     for word in query.split():
43         synsets = wn.synsets(word)
44         count = 0
45         for syn in synsets:
46             for l in syn.lemmas():
47                 if l.name().lower() not in temp.lower() and count !=2 and "_" not
in l.name() and l.name().lower() not in words:
48                     temp = temp + " " + l.name()
49                     count += 1
50     synQueries.append(temp)
51 synQueries = remove_pun(synQueries)
52 synQueries = combine_texts(improvedQueries,synQueries)
53 write_queries(originalQueries,output+"extra-20-syn-queries",synQueries,weight)
54 os.chdir(currPath)

```

Αφού δημιουργηθούν τα queries χρησιμοποιούμε την μηχανή του indri lemur ώστε να ξανακάνουμε την αναζήτηση των information need μας και συγκρίνουμε τα αποτελέσματα που θα εξάγουμε με το **trec_eval**. Στις παρακάτω εικόνες φαίνεται και το evaluation των καινούριων queries (4 σε αριθμό). Η 1^η εικόνα δείχνει το evaluation στα βελτιωμένα queries (αριστερά: έξτρα όροι, δεξιά: έξτρα όροι + συνώνυμα) έχοντας ως βάση τα αποτελέσματα των queries που είχαν μόνο τον τίτλο ενώ η 2^η είναι με βάση τα αποτελέσματα των τίτλων και των περιγραφών (έχοντας ως βάση εννοούμε ότι χρησιμοποιήθηκαν εκείνα τα αποτελέσματα για την εφαρμογή της ανάδρασης σχετικότητας).

After we have created the improved queries we use the indri search engine to rerun our information need and then compare our results with **trec_eval**. The below images represent the evaluation for each method in the corresponding field. The first image shows the evaluation of the improved queries (**left**: extra terms, **right**: extra terms and synonyms) based only on the "titles-only" queries while on the second image the improved queries are based on the "titles-and-desc" queries.

***extra**: most frequent

eval_titles_extra_20.txt			eval_titles_extra_20_syn.txt		
1	runid	all indri	1	runid	all indri
2	num_q	all 150	2	num_q	all 150
3	num_ret	all 150000	3	num_ret	all 150000
4	num_rel	all 14013	4	num_rel	all 14013
5	num_rel_ret	all 7009	5	num_rel_ret	all 7400
6	map	all 0.2376	6	map	all 0.2309
7	gm_map	all 0.1139	7	gm_map	all 0.1118
8	Rprec	all 0.2776	8	Rprec	all 0.2745
9	bpref	all 0.2521	9	bpref	all 0.2477
10	recip_rank	all 0.6258	10	recip_rank	all 0.6437
11	iprec_at_recall_0.00	all 0.6690	11	iprec_at_recall_0.00	all 0.6934
12	iprec_at_recall_0.10	all 0.4099	12	iprec_at_recall_0.10	all 0.5006
13	iprec_at_recall_0.20	all 0.3946	13	iprec_at_recall_0.20	all 0.3886
14	iprec_at_recall_0.30	all 0.3299	14	iprec_at_recall_0.30	all 0.3214
15	iprec_at_recall_0.40	all 0.2722	15	iprec_at_recall_0.40	all 0.2579
16	iprec_at_recall_0.50	all 0.2267	16	iprec_at_recall_0.50	all 0.2081
17	iprec_at_recall_0.60	all 0.1735	17	iprec_at_recall_0.60	all 0.1593
18	iprec_at_recall_0.70	all 0.1285	18	iprec_at_recall_0.70	all 0.1156
19	iprec_at_recall_0.80	all 0.0885	19	iprec_at_recall_0.80	all 0.0719
20	iprec_at_recall_0.90	all 0.0465	20	iprec_at_recall_0.90	all 0.0489
21	iprec_at_recall_1.00	all 0.0235	21	iprec_at_recall_1.00	all 0.0213
22	P_5	all 0.4493	22	P_5	all 0.4680
23	P_10	all 0.4220	23	P_10	all 0.4293
24	P_15	all 0.3987	24	P_15	all 0.4062
25	P_20	all 0.3797	25	P_20	all 0.3775
26	P_30	all 0.3351	26	P_30	all 0.3409
27	P_100	all 0.2167	27	P_100	all 0.2197
28	P_200	all 0.1516	28	P_200	all 0.1443
29	P_500	all 0.0878	29	P_500	all 0.0829
30	P_1000	all 0.0533	30	P_1000	all 0.0499

eval_titles_and_desc_extra_20.txt			eval_titles_and_desc_extra_20_syn.txt		
1	runid	all indri	1	runid	all indri
2	num_q	all 150	2	num_q	all 150
3	num_ret	all 150000	3	num_ret	all 150000
4	num_rel	all 14013	4	num_rel	all 14013
5	num_rel_ret	all 6991	5	num_rel_ret	all 6821
6	map	all 0.2084	6	map	all 0.2151
7	gm_map	all 0.0724	7	gm_map	all 0.1125
8	Rprec	all 0.2371	8	Rprec	all 0.2668
9	bpref	all 0.2259	9	bpref	all 0.2362
10	recip_rank	all 0.6002	10	recip_rank	all 0.6510
11	iprec_at_recall_0.00	all 0.6313	11	iprec_at_recall_0.00	all 0.6925
12	iprec_at_recall_0.10	all 0.4497	12	iprec_at_recall_0.10	all 0.4905
13	iprec_at_recall_0.20	all 0.3463	13	iprec_at_recall_0.20	all 0.3837
14	iprec_at_recall_0.30	all 0.2706	14	iprec_at_recall_0.30	all 0.3071
15	iprec_at_recall_0.40	all 0.2316	15	iprec_at_recall_0.40	all 0.2482
16	iprec_at_recall_0.50	all 0.1759	16	iprec_at_recall_0.50	all 0.1874
17	iprec_at_recall_0.60	all 0.1273	17	iprec_at_recall_0.60	all 0.1339
18	iprec_at_recall_0.70	all 0.0886	18	iprec_at_recall_0.70	all 0.0871
19	iprec_at_recall_0.80	all 0.0473	19	iprec_at_recall_0.80	all 0.0429
20	iprec_at_recall_0.90	all 0.0212	20	iprec_at_recall_0.90	all 0.0205
21	iprec_at_recall_1.00	all 0.0085	21	iprec_at_recall_1.00	all 0.0086
22	P_5	all 0.4507	22	P_5	all 0.4773
23	P_10	all 0.3980	23	P_10	all 0.4347
24	P_15	all 0.3636	24	P_15	all 0.4133
25	P_20	all 0.3353	25	P_20	all 0.3817
26	P_30	all 0.2996	26	P_30	all 0.3362
27	P_100	all 0.1919	27	P_100	all 0.2037
28	P_200	all 0.1360	28	P_200	all 0.1307
29	P_500	all 0.0764	29	P_500	all 0.0768
30	P_1000	all 0.0466	30	P_1000	all 0.0455

In the following matrix we see the metrics we are interested in:

Metric	titles + extra20	titles + extra20 + syn	titles + desc + extra20	titles + desc + extra20 + syn
Mean Average Precision	23.76 %	23.09 %	20.04 %	21.51 %
R-Precision	27.76 %	27.45 %	23.71 %	26.68 %
Precision@10	42.2 %	42.93 %	39.8 %	43.47 %

Based on the results above we can come to some conclusions. The extra terms (most frequent ones) and the synonyms applied to the "titles-only" original queries improved by a small percentage the MAP while the "titles-and-desc" decreased it by a lot. This can be explained by the length of our query which in the latter case is extremely large. In the R-Precision we observe the same things as the MAP. Moreover, on the "titles-only" with teh extra terms we see our best score. Finally, for the Precision@10 (which was explained in the part A of this project), we once again see similar percentages (between 40-43.5). As we can understand our improved queries, which used for the Relevance Feedback only the titles (and not the desc too), improved our MAP and R-Precision percentages and kept Precision@10 to the same levels (realistically we had a slight drop). On the other hand by using both the titles and the desc for the Relevance Feedback we dropped our metrics.