

DMT - Homework 1

Tran Luong Bang, Juan Mata Naranjo
Master in Data Science

April 13, 2021

	Cranfield Dataset	Time Dataset
Num Indexed Docs	1400	423
Num Queries	225	83
Num Queries in GT	110	80

Table 1: Overview Table

Conf ID	Text Analyzer	Scoring Functions
1	RegexTokenizer() LowercaseFilter() StopFilter(stoplist = STOP_WORDS)	scoring.MultiWeighting(scoring.Frequency(), title=scoring.BM25F(), content=scoring.TF_IDF())
2	RegexTokenizer() StopFilter(stoplist = STOP_WORDS) LowercaseFilter() StemFilter()	scoring.BM25F(K1=1.2, B=.75)
3	StemmingAnalyzer(stoplist=STOP_WORDS)	scoring.TF_IDF()
4	FancyAnalyzer()	scoring.BM25F(K1=2, B=.8)
5	SimpleAnalyzer()	scoring.TF_IDF()
6	StandardAnalyzer()	scoring.Frequency()
7	FancyAnalyzer()	scoring.MultiWeighting(scoring.Frequency(), title=scoring.TF_IDF(), content=scoring.BM25F())
8	StemmingAnalyzer(stoplist=STOP_WORDS)	scoring.BM25F()
9	RegexTokenizer() StopFilter(stoplist = STOP_WORDS) LowercaseFilter() StemFilter()	scoring.MultiWeighting(scoring.Frequency(), title=scoring.TF_IDF(), content=scoring.BM25F(K1=2, B=.8))
10	SimpleAnalyzer()	scoring.BM25F(K1=1.2, B=.75)
11	StemmingAnalyzer(stoplist=STOP_WORDS)	scoring.BM25F(K1=2, B=.8)
12	FancyAnalyzer()	scoring.TF_IDF()

Table 2: Configuration Overview

1 Part 1.1

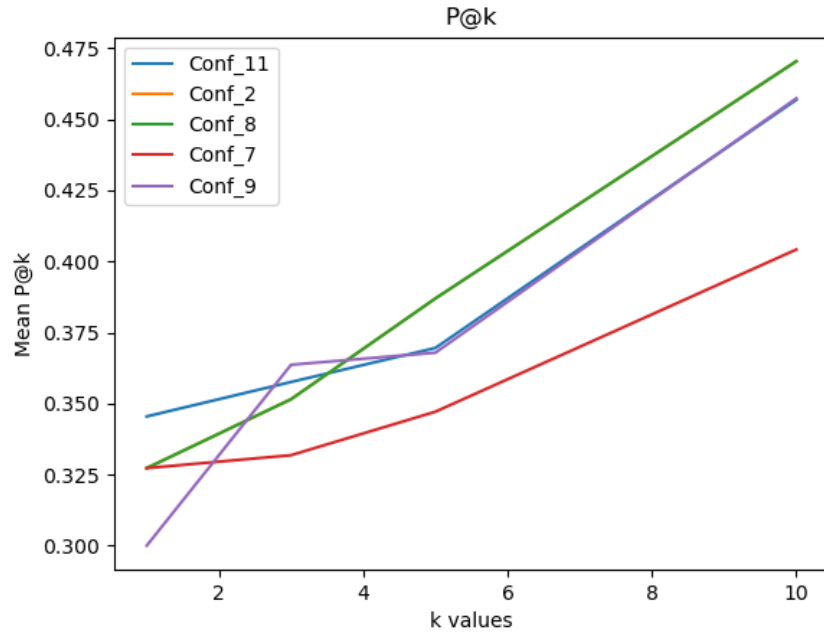
Before we deep dive into each of the individual datasets and their respective results we will start out by first giving a brief overview of the number of documents we have, the number of queries applied over these documents and finally the number of queries for which we have their ground truth at our disposal. We will see that the number of ground truth results we have are much less than the total number of queries. For the evaluation metrics we will assume that the queries for which we don't have ground truth don't exist (we remove them from all evaluation metrics). We will then also outline the different text-analyzer + scoring-functions we have deployed.

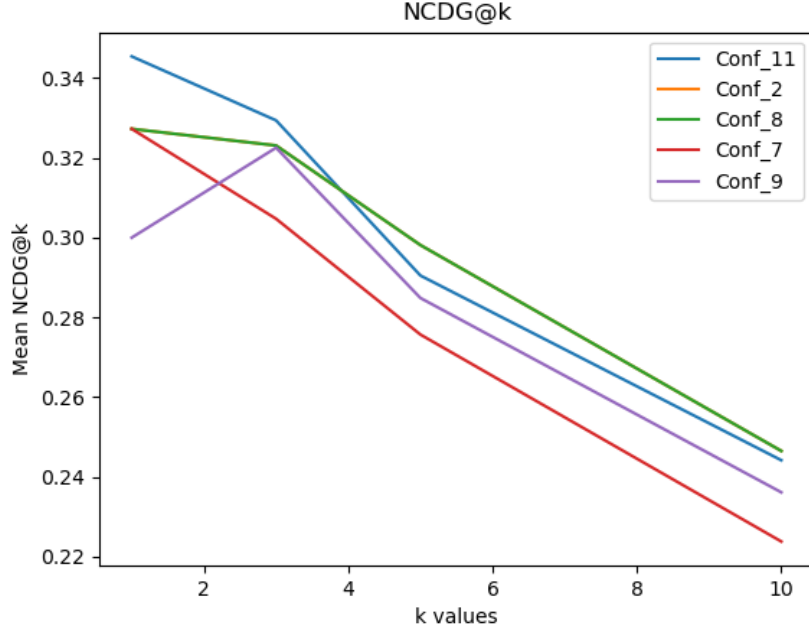
1.1 Cranfield Dataset:

We will now continue by presenting the most relevant results on the Cranfield Dataset. The first result presented is a table containing all the configurations (ordered by the MRR evaluation metric from highest to lowest), and the R-Precision metrics:

Conf ID	MRR	Mean	Min	1st quartile	Median	3rd quartile	Max
11	0.527	0.273	0	0	0.250	0.429	1.000
2	0.522	0.278	0	0	0.250	0.490	1.000
8	0.522	0.278	0	0	0.250	0.490	1.000
7	0.504	0.255	0	0	0.250	0.460	0.667
9	0.499	0.260	0	0	0.250	0.500	1.000
4	0.496	0.264	0	0	0.250	0.448	1.000
10	0.475	0.245	0	0	0.250	0.429	0.667
3	0.416	0.178	0	0	0.143	0.296	1.000
12	0.388	0.177	0	0	0.143	0.286	1.000
1	0.384	0.184	0	0	0.143	0.286	1.000
6	0.303	0.131	0	0	0.000	0.217	1.000
5	0.166	0.074	0	0	0.000	0.103	0.833

Table 3: MRR and R-Precision Overview





1.2 Time Dataset:

2 Part 2.1

The first thing we have to do before running the LSH algorithm is to figure out the best r (number of rows in each band), b (number of bands) and n (number of hash functions used) parameters such that we can minimize the number of False Positives and False Negatives. In particular, the choice of parameters will mainly help us to minimize the quantity of False Negatives. Taking into account the other constraints which are the following:

$$r \cdot b = n \quad (1)$$

$$0.97 > 1 - (1 - 0.95^r)^b \quad (2)$$

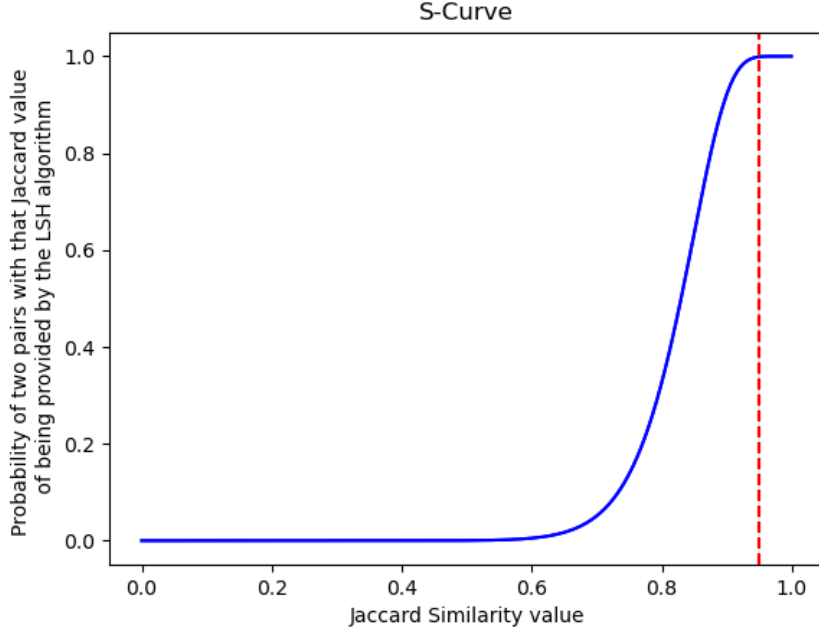
From the second constraint we can see that the following relation holds:

$$b > \frac{\log(1 - 0.97)}{\log(1 - 0.95^r)} \quad (3)$$

In order to minimize the False Negatives we will stress the 0.97 bound as much as possible, specifically up to 0.999. This will reduce the number of False Negatives, but will of course increase the number of False Positives. Since we will be able to remove some of the False Positives later on we will not worry too much for the moment. With this assumption, fixing a value of $r = 15$ we have decided to use the following parameters:

r	15
b	12
n	180

Table 4: Parameter Overview



After running the LSH algorithm with the previous set of parameters we got the following number of *potential* near duplicate matches: **34655**. However, as we already mentioned before, we can still reduce the number of False Positives a posteriori by computing the Jaccard Similarity over the set of pairs provided by the LSH tool, and therefore eliminating those pairs which we will not consider as near duplicates. For this purpose we have used the code named "False Positive Reduction" (in which we re-compute the real Jaccard similarity and not the approximated one), after which the number of *final* near duplicates was: **33868**. Of course, in this specific case we could reduce the number of False Negatives, making the LSH algorithm return more potential pairs (approx. 13%). This might be too many potential pairs in other applications where instead of 250K documents we have billion documents. In these cases the percentage of potential pairs in output must be much lower, and therefore we might have to deal with having more False Negatives.

The total time required to run the LSH tool with the parameters highlighted previously was of **7 minutes and 15 seconds** on our machine.