

11791 Design and Engineering of Intelligent Information Systems Fall 2013 Assignment 4

Kuo Liu
kuol@andrew.cmu.edu

October 28, 2013

Task 1

After dropping stop-words and changing all the words to the lowercase, I get the following result.
The output for task1 is as follows:

Score:	0.6123724356957945	rank=1	rel=1	qid=1	Classical music may never be the most popular music
Score:	0.4629100498862757	rank=1	rel=1	qid=2	Climate change and energy use are two sides of the same coin.
Score:	0.5	rank=2	rel=1	qid=3	The best mirror is an old friend
Score:	0.18257418583505536	rank=2	rel=1	qid=4	If you see a friend without a smile, give him one of yours
Score:	0.23570226039551587	rank=1	rel=1	qid=5	Old friends are best

(MRR) Mean Reciprocal Rank ::0.8

Task2

Error Analysis: Among the five test cases, two fail to rank to the relevant document to the first place. The reason is that after deleting the stop-words, the term vectors for query and the corresponding texts for case3 and case4 are:

Query3: {oneself=1, one's=1, best=1, friend=1}

Text1: {mirror=1, old=1, best=1, friend=1}

relevant

Text2: {one=1, best=2, brings=1, friend=1}

irrelevant but ranked first

Text3: {old=1, friends=1, best=1, antiques=1}

Query4: {distance=1, new=1, shortest=1, friends=1, smile=1}

Text1: {wrinkles=1, friends=1, scowl=1, smile=1, wear=2}

irrelevant but ranked first

Text2: {see=1, give=1, one=1, smile=1, without=1, friend=1}

relevant

Text3: {girls=1, put=1, behind=1, best=1, every=1, smile=1, friend=1}

For Query3, we can figure out that that the term frequency for best is 2 in Text2, while that is 1 in Text1. This point leads to the bad result. But when we turn to see the original texts, we can figure

out that old friend is the synonym of best friend (i.e. we can deem old and best as the same).

We can figure out the reason for this fail here which is we fail to consider synonyms.

So, we can add one part to do the analysis of synonyms. But since synonyms dictionaries or tools such as wordnet are not available this time, I just add one rule for simplicity. If we come up with “old friend”, we deem the word “old” as “best”. Here, we can just consider “best” as the candidate of one group of synonyms.

For Query4, we can figure out that term ‘friends’ appears in query and Text1, however, not in Text2. But after taking a deep look at the terms of Text2, we can figure out that there is a ‘friend’. So the reason for this failure is that we haven’t done morphological transformation for terms. Also, I add one rule to turn ‘friends’ into ‘friend’ since no tools to do morphological transformation are provided.

Actually, I have considered several ways to do this kind of job.

1. Use Statistical Models like BM25 or Language Models to do this kind of job. But considering that the documents are two short which results in too little statistical properties to use. I didn’t use it.
2. Use several ways to compute different similarities and sum them up, but it comes out that this way is also useless to improve MRR.
3. Use different ways to compute similarities and each way should vote for the text which gets the highest corresponding score. After the vote, we do a simple bagging. This way is also useless to improve MRR.

At last, I find out that the fundamental problem is about semantic analysis. In our case, it is synonym analysis and morphological transformation.

After these improvements, the results are:

Score:	0.6123724356957945	rank=1	rel=1	qid=1	Classical music may never be the most popular music
Score:	0.4629100498862757	rank=1	rel=1	qid=2	Climate change and energy use are two sides of the same coin.
Score:	0.6123724356957946	rank=1	rel=1	qid=3	The best mirror is an old friend
Score:	0.3651483716701107	rank=1	rel=1	qid=4	If you see a friend without a smile, give him one of yours
Score:	0.47140452079103173	rank=1	rel=1	qid=5	Old friends are best

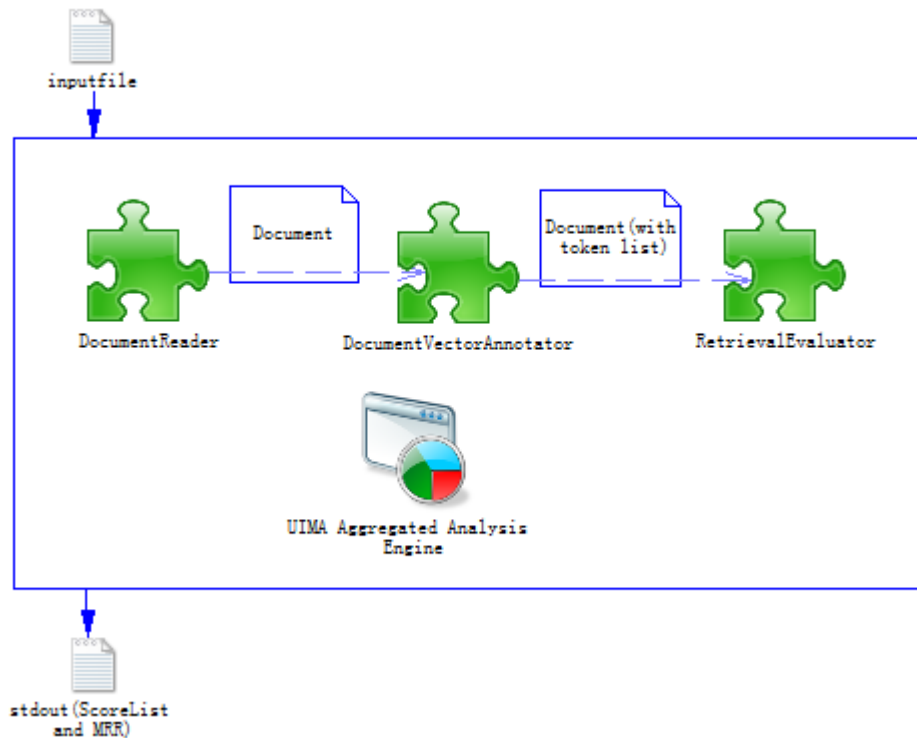
(MRR) Mean Reciprocal Rank ::1.0

Design and Implementation

The design pattern is just the same as the given template and I implemented the corresponding part and did some necessary changes.

There are two types, Document and Token, in my homework4 project.

The plot for elements and dataflow is as follows:



Bonus

I tried cosine similarity, dice coefficient and Jaccard coefficient on the test dataset using the given stop-word list and my system. The result is listed as follows (The results are about before morphological transformation and synonym analysis).

Similarity Computation	Cosine Similarity	Dice Coefficient	Jaccard Coefficient
MRR	0.80	0.90	0.90

The result of the experiment can't display the relative advantages or disadvantages among the three methods of similarity computation since the dataset is too small. Here comes the analysis for each method.

For Dice Coefficient and Jaccard Coefficient, the two methods are very similar since they both consider the degree of overlap as the measurement of similarity. The disadvantage is that they both take each element as equally important which can be a problem. Since different elements should have different weights, assigning the same weight to each element is not suitable in the usual case. Cosine Similarity corrects this disadvantage by considering the weight for each element. But it also has some disadvantages. It just considers the angle between two vectors. In this case, two vectors with the same scale are considered to be the same which can lead to problems sometimes. For example, can we deem the two vectors $\langle 1, 2 \rangle$ and $\langle 2, 4 \rangle$ as the same? In some cases, the answer may be not.