

Query Revaluation Method For Legal Information Retrieval

Liang Liu^a, Lexiao Liu^{b,*}, Zhongyuan Han^c

^a Heilongjiang Institute of Technology, Harbin, China

^b Beihang University, Beijing, China

^c Foshan University, Foshan, China

Abstract

In this paper, we introduced in detail the method of implementing the task of identifying relevant prior cases in artificial intelligence for legal assistance. For the task, we transformed the problem into a retrieval task and used the BM25 retrieval model to try to make it perform better in this task. The improved method wins second place on MAP and the second place on BPREF.

Keywords

Legal Information Retrieval, Language Model, BM25, IDF, Identifying Relevant Prior Case

1. Introduction

It is of great importance to give prior cases for the Common Law system². A prior case (also called a precedent) is an older court case related to the current case, which discusses similar issue(s) and which can be used as a reference in the current case [1]. Therefore, legal practitioners need to find and study prior cases to study how to explain current issues in older cases.

Artificial Intelligence for Legal Assistance (AILA) is a series of shared tasks aimed at developing datasets and methods for solving a variety of legal informatics problems [2]. AILA2020³, which aims to develop an automatic system and focused on precedent and statute retrievals for a given legal scenario [5]. This year AILA will consist of two different tasks. Task 1 is the same as AILA 2019 and we will focus on task 1a in this paper.

Generally, legal information retrieval is regarded as a rank task. Last year, we proposed an improvement of BM25 and achieved excellent results [3]. As early as 2004, the multiple weighted fields base on BM25 were proposed by Robertson [4]. So we will continue to improve the method of last year in 2020.

2. Methods

For the task of Identifying relevant prior cases, we treated it as an information retrieval task and submitted three runs with BM25.

2.1. Data Pre-processing

According to the statistics, the query, which is a description of the situation in *Query_doc*, contains over 500 words on average, and the document, which the prior case in *Object_casedocs*, contains over 3,000 words on average. For traditional retrieval, the query sentence in the task is too long.

Consequently, we should preprocess the data to shorten the length of the sentence without losing its main meaning. As we all know, the common method is to remove all stop words, we also chose

Forum for Information Retrieval Evaluation 2020, December 16-20, 2020, Hyderabad, India

EMAIL: trueliuliang@gmail.com (A. 1); liulx15@yeah.net (A. 2)(*corresponding author); hanzhongyuan@gmail.com (A. 3)



© 2020 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)

² https://en.wikipedia.org/wiki/Common_law/

³ <https://sites.google.com/view/aila-2020/track-description>

this method and converted the text to lowercase. Finally, we use *Lucene toolkit*⁴ to index the document.

2.2. double_liu_2020_1

For this submission, we chose the BM25 model and improved it by modifying its relevant calculation, as follows:

$$BM25(D, Q) = \sum_{i=1}^n IDF(q_i) \cdot \frac{TF(q_i, D) \cdot (k_1 + 1)}{TF(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgdl})} \quad (1)$$

In this formula (1), we given the definition of BM25, where q_i is the word in Q , $|D|$ is the length of document D , and $avgdl$ is the average document length in the text collection, k_1 and b are the parameters of BM25. In this task, we set the parameter $k_1=2.99$ and $b=0.65$.

Furthermore, we modify the relevant computation to get an improved BM25.

$$rel(D, Q) = BM25(D, Q) + BM25(D, Q') \quad (2)$$

where Q is a query sentence with stop words removed, and Q' is a keyword that is further extracted from Q , here we choose the *IDF* algorithm to sort the words in Q , and form the top $m\%$ words into Q' . m is a free parameter, and we set $m=50$.

2.3. double_liu_2020_2

Inspired by the former, we split the method in *double_liu_2020_1* into two sub-methods as our *double_liu_2020_2* and *double_liu_2020_3*.

In the *double_liu_2020_2* submission, we chose the first half of formula (2) to form our method one, as follows:

$$rel(D, Q) = BM25(D, Q) \quad (3)$$

All the other settings are followed *double_liu_2020_1*.

2.4. double_liu_2020_3

For this submission, we choose the second half of formula (2) to form method three, as shown below:

$$rel(D, Q) = BM25(D, Q') \quad (4)$$

All the other settings are also followed *double_liu_2020_1*.

2.5. Other methods

We also tried other experiments, but the results were not satisfactory.

2.5.1 Cosine Similarity

For this method, we want to rank the cosine similarity between the query sentence and the document as an indicator. Firstly, we use the bag-of-words model to construct word vectors for the query sentence and the document respectively and then calculate the cosine similarity and rank. The formula for cosine similarity is as follows:

⁴ <https://lucene.apache.org/>

$$\text{rel}(D, Q) = \text{Cos}(A, B) = \frac{A \cdot B}{|A| \times |B|} = \frac{\sum_{i=1}^n (A_i \times B_i)}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}} \quad (5)$$

Where A and B are two vectors.

2.5.2 Generalized Jaccard Similarity

In this method, we choose to use generalized Jaccard similarity as an indicator to sort.

$$\text{rel}(D, Q) = J(A, B) = \frac{\sum_{i=1}^n \min(A_i, B_i)}{\sum_{i=1}^n \max(A_i, B_i)} \quad (6)$$

2.5.3 Cosine with Jaccard

In this method, we improve the previous two methods and introduce the parameter k. The specific formula is as follows:

$$\text{rel}(D, Q) = k \times \text{Cos}(A, B) + (1 - k) \times J(A, B) \quad (7)$$

where k is a free parameter, and we set k=0.3

3. Results

3.1. Evaluation Measures

Standard Information retrieval metrics like Measures like Precision, Recall, Mean Average Precision (MAP)⁵, Discounted Cumulative Gain(DCG) and Mean Reciprocal Rank(MRR) will be used for evaluation in the task.

3.2. Evaluation Results

Table 1. Results of the AILA Task 1a sorted by MAP

| Run_ID | MAP | BPREF | recip_rank | P@10 | rank |
|-------------------|---------------|---------------|---------------|-------------|------|
| double_liu_2020_3 | 0.1382 | 0.1045 | 0.1886 | 0.07 | 2 |
| double_liu_2020_1 | 0.1306 | 0.0737 | 0.1963 | 0.07 | 4 |
| double_liu_2020_2 | 0.123 | 0.0621 | 0.1969 | 0.08 | 11 |
| Jaccard | 0.0820 | 0.0578 | 0.1464 | 0.07 | - |
| Cosine_Jaccard_k | 0.0781 | 0.0563 | 0.1521 | 0.08 | - |
| Cosine | 0.0490 | 0.0116 | 0.1579 | 0.04 | - |

⁵ <https://trec.nist.gov/pubs/trec16/appendices/measures.pdf>

Table 2. Results of the AILA Task 1a sorted by BPREF

| Run_ID | MAP | BPREF | recip_rank | P@10 | rank |
|-------------------|---------------|---------------|---------------|-------------|------|
| double_liu_2020_3 | 0.1382 | 0.1045 | 0.1886 | 0.07 | 2 |
| double_liu_2020_1 | 0.1306 | 0.0737 | 0.1963 | 0.07 | 9 |
| double_liu_2020_2 | 0.123 | 0.0621 | 0.1969 | 0.08 | 16 |
| Jaccard | 0.0820 | 0.0578 | 0.1464 | 0.07 | - |
| Cosine_Jaccard_k | 0.0781 | 0.0563 | 0.1521 | 0.08 | - |
| Cosine | 0.0490 | 0.0116 | 0.1579 | 0.04 | - |

It can be seen from Table 1 and Table 2 that double_liu_2020_3 is the best among the methods we submitted. According to the results, the relevant prior case information is helpful to guide the judgment of current case.

4. Conclusion

In this task, we describe a method that uses an improved BM25 to identify relevant priors, and it can be concluded that using certain algorithms to extract keywords will improve efficiency. Compared with other submissions of the task, our improved BM25 model can get the second place in MAP and BPREF.

5. Acknowledgments

This work is supported by National Social Science Fund of China (No.18BYY125).

6. References

- [1] Mandal, A., Ghosh, K., Bhattacharya, A., Pal, A., Ghosh, S.: Overview of the fire 2017 ired track: Information retrieval from legal documents//Proceedings of FIRE 2017 - Forum for Information Retrieval Evaluation, 2017:63-68
- [2] Bhattacharya, P., Ghosh, K., Ghosh, S., Pal, A., Mehta, P., Bhattacharya, A., Majumder P.: Overview of the FIRE 2019 AILA Track: Artificial Intelligence for Legal Assistance//Proceedings of FIRE 2019 - Forum for Information Retrieval Evaluation, 2019.
- [3] Zhao, Z., Ning, H., Huang, C., Kong, L., Han, Y., Han, Z.: Fire2019@aila: Legal information retrieval using improved bm25//Proceedings of FIRE 2019 - Forum for Information Retrieval Evaluation, 2019:40-45.
- [4] Robertson S, Zaragoza H, Taylor M. Simple BM25 extension to multiple weighted fields//Proceedings of the thirteenth ACM international conference on Information and knowledge management. 2004: 42-49.
- [5] P. Bhattacharya, P. Mehta, K. Ghosh, S. Ghosh, A. Pal, A. Bhattacharya., P. Majumder, Overview of the Fire 2020 AILA track: Artificial Intelligence for Legal Assistance. In Proc. of FIRE 2020 - Forum for Information Retrieval Evaluation, Hyderabad, India, December 16-20, 2020.