

Evaluating pre-trained transformer models for similar article recommendation in PubMed

Li Zhang

May 22, 2021

Abstract

The work presented here promises to make the exploration of scholarly material faster and more accurate. Though the simialr method identifcaiton method shows in this paper is for PubMed, we believe it can be migrate into other bibliographic databases/digital libraryies and border field.

1 Introduction

Similar article recommendation is an important feature in many academic searching database/digital libraries. It enable users to go through relative researcher quickly, play an important role in improving user searching experience, while more importantly, it can help to accurately dismiss the valuable biomedical findings. Improved literature search engines can save researchers time and effort by making it easier to locate the most important and relevant literature. [?]

How long has PubMed supplied "Similar article" function? Like other system, ResearchGate, Scopus, NCBI's PubMed system, has integrated this feature since *. However, their method to find the similar article is still very *. To date, with the fast development of natural language processing, the cutting-edge techniques have provided an opportunity to improve similar article recommendation performance.

In this paper, we show the similar article performance can be largely promoted by using large-heavy pretrained language model.

PubMed related article links identify closely related articles and enhance our ability to navigate the biomedical literature [?].

PubMed has integrated the "Similar article" function for a long time, show the case here.

significance: Fig * shows the similar article recommendation functionality, this function is very helpful for biomedical scholar, as a recent works by NLM/NCBI team suggest user needs can be largely improved while users explore related articles. can power PubMed user experience,...

Why similar article recommendation is very important for further improving search experience, why it is a necessary functionality? Related works: how PubMed improving user searching experience. To improving user searching experience, NCBI has provided many in place measurements from several aspects. such as auto-completion, ..., frequent search terms recommendation. However, these measures can be not necessary enough, user may explore other What did European PubMed did? and what did other platform did?

In many academic service platform, such as Web of Science, ..., they How did others find similar article?

Our contribution are three parts. we provide a method to automatically build similar article dataset for development models we evaluated the most well-known pre-trained models on four dataset, and empirical evaluation shows fine-tuned * model shows state-of-the-art result. Using this method, we obtain the paper distribution vector for whole PubMed papers.

This paper targets for three groups of readers. 1. NCBI's PubMed, we created an enhanced similar article methods for PubMed, and show better results than the same current result in PubMed. Moreover, we successfully compress the original 768-long vector into a short concise vector for entire PubMed paper. This exercise helps to make this method practical. Based on the concise vector, we also pre-computed the "nearest" paper for all PubMed papers, and release the paper vectors for facilitate future studies. 2. Research group for science of science. This paper helps them to better visualize the clustering of indexed paper in PubMed, and help them to manage/maintenance. For each paper, one can use the tool in our released code to fast identify the similar articles by similarity in descending order.

In the following section. We review existing or potential methods for this task, and review evaluation datasets in Related work section. Then, in method section, we show how to building distribution representation for articles, and how to build our evaluation dataset.

2 Related works

2.1 Method review

2.1.1 PubMed Similar Article Function

How NCBI did? To identify similar article in PubMed, the NCBI group has invested many efforts. According to the recent description of the implementation of this functionality^{1, 2}, measurement like word-level * is used to compute the similarity. However, this solution may be not an optimal one, because the words extracted from citation abstract already helpful for accurately similarity because some general words can convey the research content appropriately. especially for synonyms/, it will result in the fact that similar paper will be removed due to words with close meanings appear in similar articles. They mainly rely on lexical overlap, they cannot learn the semantic similarity, this is especially occurs for some complexity compound and long biomedical phrase.

How other researchers did for PubMed? In recent years, there are several studies on similar paper identification task for PubMed.

Assume this task is text similarity problem, is there any methods can do it better? From a broader view, this task is a specific case of text similarity problem, i.e., text similarity for biomedical articles. Therefore, techniques

¹<https://pubmed.ncbi.nlm.nih.gov/help/computation-of-similar-articles>

²<https://pubmed.ncbi.nlm.nih.gov/help/similar-articles>

on text similarity can be borrowed for improving the experience of similar article exporelation in PubMed. In addition to the existing studies for this task, Recent advances on natural language processing has achieve great progress, many effective models such as deep text matching, Bert cite are proving to be effective on similar text task. These models have removed the barrie of words mearning variants and close words meaning between articles by considering the simialrity from a semantic prospective. What we can do for improving it? In this work, we investigate the cutting-edge text similarity approaches in NLP for biomedical arcicle and evaluating them for this helpful functionality.

2.2 Dataset review

Show detailed is there any evaluation dataset? how others bulding the dataset? similar paper finding as a information retrivel problem.

3 Method

We employ SBert model ³, a increasing number of state-of-the-art pretrained models for more than 100 languages, fine-tuned for various use-cases. Specifically, BioBert, SciBert and ... are introduced for comparsion. Since these models are pre-trained on a vast scientific-text corpus for general tasks. it may not achieve perfect performance on the specific task. Therefore, in order to obtain better results, we fine-tune these models and evlauate the original models and the fine tuned models seperatly.

3.1 dataset

We ranomdly selected 1M articles from PubMed citations, and developed programs to bulk retrieve the similar articles for the 1M samples from PubMed official site using its open APIs via biopython [?] package.

In this section, we show how we did to bulid the dataset. For evaluation dataset. We ...

Note that we did not consider the publication timeframe, as we can a later-published article can exist in the related article list.

plot a database enhanced here by NER and Euopen PubMed.

Rules to selected ground truth simiar paper for buliding evaluation dataset. In order to buliding evaluation dataset. We used two rules to identify the

³<https://www.sbert.net/>

similar paper from entire PubMed citations. These rules are inspired by the characteristics of similar papers.

Actually, very similar paper usually share similar research purpose, and they use similar methods or tools. These identifiable features can be extracted by the representative terms-MeSH, which is a standard resource for the whole spectrum of biomedical scientific papers. Currently, the MeSH terms in each citation are assigned carefully by the NLM's curator team. MeSH are used for widely applications not restricted to paper search. MeSH are indexed through the collaboration and machine-aided indexing and human curators. Another prominent advantage of MeSH is semantic indexing over lexical, the same MeSH term can be labeled to a paper though it does not exist in the paper's title and abstract.

In addition, the citing-cited relationship between articles is another feature needing exploration for building dataset. It can be understood that the cited articles in reference usually have semantic connection for the citing paper. Thus, the relationship can be much helpful for building and accurate paper similarity dataset.

However, neither MeSH list or reference list is very sparse for PubMed paper. On average, there are * percentage paper having MeSH, and there are only * number of MeSH for these articles. For citations, the low density of citation relationship of PubMed is a well known problem. Similar to MeSH list, only * percentage paper having reference, and only average * number of citations for these articles.

To handle this issue, we made considerable efforts to enhance the two metadata. For MeSH enhancement, we use the NER result of PubMed Knowledge Graph [?]. The goal of the NER task is to recognize named entity from PubMed citation content (i.e., abstract). Our empirical evaluation found there are nearly *% entities included by MeSH. To supplement the additional MeSH, we first filter MeSH terms from NER entity list and then associate them to the whole PubMed citations by the PubMed paper ID (PMID). With respect to the reference sparsity problem, we compensate it by retrieving the reference information from Europe PMC⁴, which is a service of the Europe PMC Funders' Group, in partnership with the European Bioinformatics Institute; and in cooperation with the National Center for Biotechnology Information at the U.S. National Library of Medicine (NCBI/NLM). It includes content provided to the PMC International archive by participating publishers. By the open API⁵ of Europe

⁴<https://europepmc.org/>

⁵<https://europepmc.org/RestfulWebService!/Europe32PMC32Articles32RESTful32API/references>

PMC, We are able to obtain a relatively complete reference list. Through this exercise, the overall density of reference for PubMed is * in comparison to *% of original reference density.

To incorporate it. we

In this paper, Enhanced Mesh and Enhanced Citation. Simialr paper usually has.

However,

MeSH/Citation

3.2 fine-tuning

4 Results

4.0.1 Dataset

4.0.2 Competing methods

As discussed in relative words section, there are many methods in NLP for text similarity. For evaluation, we introduce Bert, For comparison, we introduce the recommendation result shown in PubMed interface (see "Simialr article" in page navigation of a particular paper)

4.0.3 Evaluation metrics

Ranking metrics, Information retrieval metrics and paper title overlapping

5 Discussion

5.1 how this works can be integrated into PubMed system?

fast, very short words embeddings, inductive-infering.

5.2 user study

how it can be

5.3 limitation

5.4 future potential improvement

This study only recommendation paper from semantic perspective, however, in many commercial recommender system, the recommended items may not

only semantic relatedness, other relationship the relationships that a mature system should consider is not equal to semantic relatedness. other crucial aspect such as ... also need consideration. However, we can not obtain such real word dataset, i.e, integrating PubMed searching log to develop more powerful recommender system. Thus, we can image the paper recommendation system in PubMed can be more powerful by leveraging state-of-the-art technique in recommendation system and information retrieval. for example, recommendation with interpretation,

exploring more user intelligence that can be available in NCBI

6 Conclusions

In this works, we shows an effective and efficient way to locate similar PubMed for power user search experience in PubMed system. This study provide initial, portary investigation on simiar article for PubMed system. ... Our future works intent more relationship in paper recommendation.