

CSC 510 Software Engineering: Project 1 ^{*}

Group-P [†]

Monica Metro North
Carolina State University
1932 Wallamaloo Lane
Wallamaloo, New Zealand
mgmetro@ncsu.edu

Zachery DeLong North
Carolina State University
2305 Horizon Hike Ct
Raleigh, NC
zpdelong@ncsu.edu

Zhangqi Zha North Carolina
State University
1932 Wallamaloo Lane
Wallamaloo, New Zealand
zzha@ncsu.edu

Bikram
Brookhaven Laboratories
Brookhaven National Lab
P.O. Box 5000
bikram@ncsu.edu

ABSTRACT

Choosing a textbook for a class is something that many professors do several times a year, but the vast number of available textbooks for any given subject makes choosing one that both meets the needs of the class and meets the budget constraints of today's students is a challenge. In this paper, we propose two methods for automatically learning topics from books, a set of evaluations for those methods, and a system that will use those methods in a web app that attempts to suggest textbooks to a professor that balance hitting all topics the course requires while still being economical for students.

1. INTRODUCTION

One of the most basic tasks that a professor must do is identify what textbook to use for a given class. While some fields have textbooks that have become popular and are clearly the best in their subject matter, many fields (especially emerging ones) have no such exemplar. It would be a complex enough problem if there were not already a massive number of textbooks on sources such as Amazon, but self aware professors have attempted to use textbooks that are less expensive with hopes of allowing lower income students to attend more easily. This admirable intention serves to increase the time needed to research the given books, and there is no obvious heuristic to apply to searches to limit the field.

^{*}(Paper 1). For use with SIG-ALTERNATE.CLS. Supported by ACM.

[†]A full version of this paper is available as *Author's Guide to Preparing ACM SIG Proceedings Using L^AT_EX2_ε and BibT_EX* at www.acm.org/eaddress.htm

In the interests of making students lives less expensive, we propose exploring a system which can identify topics in books automatically, and which, given a set of requirements from a user (usually a professor), can suggest options of varying expense and completeness for a given set of topics. To do this, we propose exploring word clusters generated by Doc2Vec, a family of common natural language processing (NLP) algorithms that attempt to cluster similar words, and topics generated by latent dirichlet allocation (LDA), another common NLP algorithm that directly attempts to identify topics in text, to automatically infer the content of a set of textbooks. We then intend to build a web app that will allow a user to specify a set of topics and which will use the mentioned algorithms to make suggestions while minimizing the cost of said textbooks.

2. LITERATURE REVIEW

We propose two methods for analyzing and storing topics from books automatically: LDA and Doc2vec. We also give some preview of similar work that has already been done.

2.1 LDA

2.2 Doc2Vec

2.3 Related Work

3. PROPOSED SOLUTION

To tackle this problem, we propose a three-part system that will go about recommending books via a web app. The first component of this system is a recommendation engine, which will implement the above algorithms in some way and store their results (the topics) in an indexed database which we can summarily search. The next component will be a recommendation engine which will implement our recommendation algorithm (see below for more details). The third component will be a web app that will allow the end-user to search the database of topics and organize them into a useful search of books. It will then present a list of suggested books using the infrastructure mentioned earlier.

3.1 Book parser

The first major component of the system is the book parser. We intend to implement this taking advantage of the pre-built libraries available in Python such as SKLearn, Pandas, and others. The goal of the book parser is to implement one of the two algorithms discussed in section 2 and to store the parsed topics in a database for searching later. This will be implemented using the command pattern, which provides a simple interface to implementing different algorithms in code and allows us to easily swap between implementations.

The books being parsed will come from a freely available online database of textbooks. It is also our intention to cache the topics parsed from these algorithms somehow (either in a database or in a flat file somewhere) to be referenced in the future. This will allow us to train the algorithms on the books ahead of time and will make searching for books based on topics much more efficient.

3.2 Recommendation Engine

The next step in our problem is to use the topic analysis data outlined in the previous section (3.1) to suggest books based on the topics covered and the cost. This algorithm will be written in Python as well.

Conceptually, this algorithm will need to take in a set of required topics and weighting for topics vs price. This value will be used to favor inexpensive books over complete topic coverage in our algorithm. It will then search its algorithms for books that meet the required topics. The algorithm should then rank the books using the weighting mentioned earlier, then it will sort some number of the top results and give them back to the UI.

3.3 Web App

The web app is the final component of this architecture and it serves to tie the previous two components together. The app provides the UI to the algorithms, allowing professors to log in and select a set of subjects they wish to teach on, then it should present them with a UI for searching for topics in our database. It should then allow them to create a list of topics to suggest books for. Once the list has been generated, the tool should search against the book databases it has access to using the algorithms outlined in the previous sections.

We will develop this app in either the MEAN stack or some variant thereof. (Does it make more sense to just use Django since we are writing learning algorithms in Python?) The front-end will be developed using Bootstrap to save time in developing a useful UI.

4. EVALUATION PLAN

This is a test of our evaluation plan section

4.1 Evaluating Topic Modeling

4.2 Evaluating textbook suggestions

5. CONCLUSION

5.1 Anticipated challenges

5.2 Tentative schedule

5.3 Future enhancement

One major enhancement would be the introduction of a rating system to the suggestion engine. Another major enhancement would be including links in the UI to buy the books. A third major enhancement would be evaluating other topic modeling algorithms to see how they stack up. Another enhancement would be to test an even more rudimentary topic modeling method, such as simply representing topics as lines in a table of contents. Essentially asking the question of whether or not fancy algorithms are even more useful than simple text search for this domain.