

Automated Reading Comprehension Clustering

Charles Hathaway, David Hedin

December 15, 2015

Abstract

Determining the level of readability of documents, especially books, has lots of application in the domain of education. It helps to quantify and group books which may be used at a particular reading level, therefore enhancing the classroom experience for both instructors and teachers. In this paper, we ran an experiment with a variety of books obtained from Project Gutenberg [1] organized by the project maintainers into a 2 groups; books for children, and adult fiction. To further enhance the analysis of this project, we also used our features to try and cluster books into clusters defined by the FleschKincaid readability metrics [2].

1 Introduction

There are several systems currently used to classify books based on reading comprehension level, for numerous applications ranging from selecting books for classrooms, to measuring an individuals literacy skills for both medical (autism, dyslexia, etc.) and educational purposes. In this paper, we analyze the results of utilizing several existing features to classify works in addition to a variety of novel features we created. The primary goal is to cluster books in groups representing the original designation of books; adult fiction and children fiction.

Given an input of 2292 books (320 children, 2002 adult, with 32 books overlapping) we achieved an F-score of 87.5%; this is a significant improvement over the 62% baseline ¹

2 Previous Works

Although not extensively, we did evaluate and learn from a number of previous works. Most significantly, we borrowed features from work done by Feng et al. [3]. Feng focused on the concept of documents being difficult to read due to

¹There was some in-class discussion which suggested our baseline would be (total number of adult books)/(total number of books), which would put the baseline at around 86%. After testing this experimentally, and reasoning things out, we concluded the true baseline would be 62% as our algorithm had no idea what the sizes of the clusters were, and a truly random distribution would but half in each cluster, with one cluster having a higher chance of being correct than the other

”items fall[ing] out of memory before they can be semantically encoded”. With this in mind, most of the novel metrics they designed focus on the number of entities mentioned in a document. Their work was motivated to help organize documents by readability for adults with intellectual disabilities.

Our work is different than this in that we are using a much different data source (they were focused on news articles, where we are focused on fiction books). In addition, instead of focusing on adults with intellectual disabilities directly, we are instead more interested in examining books in relation to each other, with the assumption that easier to read books would be easier for both adults with ID, and developmentally delayed children.

3 Methodology and system design

The system is written in Python, and utilizes the Natural Language ToolKit (NLTK) [4]. The toolkit was trained using a portion of the Penn Treebank and the Conll2000 corpus. The portion used was provided by NLTK via it’s distribution manager (nltk.download).

Once the corpus was downloaded and configured, we built the system in 3 stages. The first was retrieving the document to be analyzed; this was done ahead of time to verify the availability of the resources, and allow us to develop without worrying about going over any kind of bandwidth quota. However, now that development is mostly complete, this step can be skipped and documents can be downloaded directly from a Gutenberg mirror.

The next step is to process the data. This includes using the trained POS tagger, chunker, and NE tagger on all documents, running a variety of feature function on the documents, then recording the device. Our system writes data to a CSV file as it is processed, and uses a ”Book” object to cache the results of each computation between features.

And lastly, we score the results of the system. This process involves reading in the CSV file and training our clustering algorithm, then clustering our test data. Ultimately, we determine the accuracy of the clustering using a B^3 scoring system.

To get a better idea of how each feature contributes to the final score, we run the system multiple times with every combination of features.

3.1 System Usage

To setup your system, please follow the instruction in the README.md file provided in the repository.

Once configured, simply run ”python src/main.py -help” to get a list of options and configuration settings. The simplest invocation of this command requires you to specify 2 files; the list of adult books, and the list of children books. It will then print the clustering results with the default feature-set to standard out.

4 Linguistic Features

Average number of words per sentence
Average number of syllables per word
Percentage of words with more than 3 syllables
Average number of noun phrases per sentence
Average number of common and proper nouns per sentence
Average number of verb phrases per sentence
Average number of adjectives per sentence
Average number of conjunctions per sentence
Average number of prepositional phrases per sentence
Total number of noun phrases in document
Total number of common and proper nouns in document
Total number of verb phrases in document
Total number of adjectives in document
Total number of conjunctions in document
Total number of prepositional phrases in document
Number of entity mentions in document
Number of unique entities in document
Average number of entity mentions per sentence
Average number of unique entities per sentence

Table 1: List of possible features from previous work

Average word length in document
Total number of unique words in document
Ratio of unique words to total number of words in document
Ratio of proper nouns to common nouns in document
Length of document
Average number of proper nouns per sentence
Total number of proper nouns in document
Total number of passive sentences in document
Average number of prepositional phrases per sentence
Total number of prepositional phrases in document

Table 2: List of possible new features

5 Results

Due to sheer number of results, we merely provided a synopsis consisting of the highest, middle, and worst combination of features.

f-score	recall	precision	combination of filters
0.8748384028	0.8787276342	0.8709834469	combination 1
0.6690611348	0.7238314176	0.6219964193	combination 2
0.6041246854	0.74973942	0.505873669	combination 3
0.622534193098	0.508873237259	0.801571724969	Baseline

Table 3: Results of clustering

Feature combinations:

1. average_number_of_adjectives + average_number_of_common_proper_nouns
2. total_number_of_adjectives+average_verb_phrases + total_number_of_noun_phrases
+ total_number_of_prepositional_phrases + total_number_of_verb_phrases
+ total_number_of_conjunctions+average_number_of_common_proper_nouns
3. average_verb_phrases + ratio_of_common_to_proper_nouns+average_number_of_conjunctions
+ average_prepositional_phrases

6 Conclusion

References

- [1] (Dec. 12, 2015). Project gutenber, [Online]. Available: <http://www.gutenberg.org/> (visited on 12/12/2015).
- [2] J. P. Kincaid, R. P. Fishburne Jr, R. L. Rogers, and B. S. Chissom, “Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel,” DTIC Document, Tech. Rep., 1975.
- [3] L. Feng, N. Elhadad, and M. Huenerfauth, “Cognitively motivated features for readability assessment,” in *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics*, Association for Computational Linguistics, 2009, pp. 229–237.
- [4] E. Loper and S. Bird, “Nltk: The natural language toolkit,” in *Proceedings of the ACL-02 Workshop on Effective tools and methodologies for teaching natural language processing and computational linguistics-Volume 1*, Association for Computational Linguistics, 2002, pp. 63–70.