

Automated Reading Comprehension Clustering

Charles Hathaway, David Hedin

October 30, 2015

1 Executive Summary

This project will focus on discussing and exploring the effectiveness of various language features when applied to clustering books based on reading comprehension level. It will utilize previous works in the area, which although may not be perfectly applicable, have a somewhat obvious connection to the objective. To test the features, we will attempt to generate 3 clusters from plain-text books freely available online; child-friendly, young adult, and adult books.

The timetable for this project is given at the end of the document.

2 Goal

The goals of this project are:

- Exploration of existing reading comprehension features and clustering tools
- Testing of new and existing reading comprehension features
- Analysis of results and discussion
- Interface to allow others to utilize the tool

3 Background and Motivation

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. The goal of this project is to enable a larger selection of books for these purposes by automated the system of clustering books by grade level, which is currently manually done at the cost of the publisher.

4 System Architecture and Approach

Following many other NLP applications, this project will utilize a pipeline approach. This pipeline is laid out in figure 6. The design is further described in the box chart, displayed in figure 6.

In order to accelerate development, the coding will primarily be done in Python. The project will utilize various open source toolkits, including the NLTK, which provides a maximum entropy clustering framework. This will allow more time to focus on selecting effective features, and less time remaking the wheel. A preliminary list of possible features used in previous work is available in table 1, and a list of newly generated features has been compiled into table 2. Some or all of these features may be implemented to test their effectiveness at clustering the data. It should be noted that many of these features were borrowed from the previous works done by Feng [1]

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

5 Deliverables

The deliverables of this project will have 3 key parts:

- A technical report, detailing the results of key experiments and output of several configurations of the systems

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

- A command-line interface that allows the scanning and clustering of a multitude of books
- A web interface which outputs the cluster a particular uploaded books belongs to

6 Required Resources

The most important resources for this project is the books themselves. To that end, Project Gutenberg will be mined to provide both the training and testing data sets. They have an extensive list of children’s books, which we will use as the basis for the first clustering. Ideally, we will end up with a clustering such that children’s book all fall below the average reading difficulty, at average - 1 standard deviation. For the second clustering, we will determine it by picking a range in between the children books and the highest threshold; ideally in the range of average ± 1 standard deviation. The last section, adult/advanced readers, will consist of books deemed more complicated than the average + 1 standard deviation.

It should be noted that the Project Gutenberg repository is large; upwards of 650GB. Most books in the project have multiple translations, which contributes to the large size (in addition to multiple translations, they also have multiple format). To make this project more achievable, we will limit ourselves to the fiction category; which has it’s own subcategory, children’s fiction.

Lastly, we need a toolkit which provides many resources to help calculate the clusters given our feature set. For this purpose, we will use the Python Natural Language Toolkit (NLTK), which is a free and open source library intended for these types of application.

References

- [1] 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.

Figure 1: Flowchart indicating flow of information within system

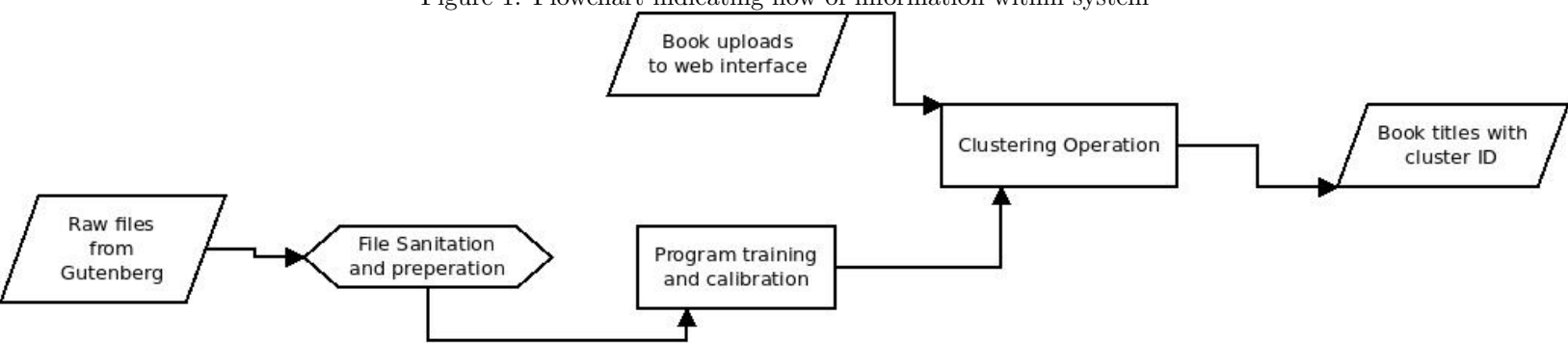


Figure 2: Box diagram of intra-program modules

