

Class project: Complex Word Identification

Pierre Finnimore

Abstract

1 Introduction

Initial words: 271

What is the task and why is it important?

The task is as follows: given a target word (or set of words) within a sentence, identify if the target is “complex”. The data given for this task is a set of labelled sentences with targets. The labels were derived from a survey of both Native and Non-native speakers of two languages: English and Spanish.

We are interested in identifying word complexity for several reasons. Automatic extraction of complex terms could help with automated tutoring systems, Natural Language Generation, writing editing software, studies into second-language acquisition, political speech analysis, Machine Translation, as well as linguistic or psychological studies into the what people find complex.

2 Baseline system description

System descriptions in enough detail for the reader to be able to understand how to reimplement your baseline models and to appreciate why they are suitable for the task at hand.

The baseline system we developed was somewhat more advanced than the initial baseline provided. Both languages used an identical model and set of features. These were:

1. Number of characters in the target
2. Number of words in the target
3. “Rarity” score based on letters of the target
4. Maximum number of consecutive vowels
5. Maximum number of consecutive consonants

6. Number of synonyms
7. Common English and Spanish Suffixes
8. Common English and Spanish Prefixes
9. Common English and Spanish Infixes
10. Common Latin and Greek prefixes

For my baseline system, I focused on features that can be extracted using only the target word; no context words were considered. First, I improved the basic architecture of the provided baseline to allow easy pattern-matching of prefixes, infixes and suffixes, and for further features of this type to be easily added.

My initial hypotheses to try for the baseline were as follows:

1. Letter rarity - this feature was chosen because words with rarer letters might be harder to understand. This feature could be seen as a vague approximation of word rarity, which would otherwise require more data to learn.
2. Max consecutive consonants, max consecutive vowels - the idea behind this was that dense combinations of letters might be difficult to parse. For example “queueing” or “rhythms”. These words deviate from the more straightforward consonant-vowel-consonant-vowel pattern.
3. Uniquely English vowel/consonant diagraphs/consonant blends - this was chosen because a non-English speaker might find these tricky. This feature could be expanded (and made more adaptable to different languages) by analysing a corpus for particularly rare combinations.

4. POS tags - because some POS are part of closed sets, and so potentially easier to understand. In addition, certain rarer tags might be harder to comprehend, especially in languages where the overall form of the word is changed, depending on its POS.
5. Number of synonyms - The idea was that ambiguous words might be more confusing. Conversely, perhaps ambiguous words are actually less likely to be regarded as complex, since it is more likely that the person knows at least one of the meanings.

While individually, these features did allow the system to predict with greater-than-random accuracy, most of them did not improve over the initial baseline. Some of the features I tried, such as Spanish prefixes, reduced the testing accuracy if they were added. This may just be noise in the data, or it may be that learning is made more difficult if irrelevant or common features are added. This may be less of a problem with a different learning model.

3 Improved system motivation and description

One of the features added for

We found a list of “1000 simple English words” on Wikipedia, and used these to extract sub-word features of two or three letters long. The idea was that these words may contain certain patterns that people regard as simple, and each sub-word feature was weighted by how often it occurred in the text. Interestingly, introducing these features improved the Spanish F-score, but not the English F-score. This may be because those combinations also happened to correlate with easy or difficult morphological features in Spanish, or it may be that the participants of the original study found that combinations that they recognised from English were easier or more difficult to identify in Spanish.

The single largest improvement to the model was obtained by using a Random Forest Classifier rather than a Logistic Regression classifier. There are several reasons that we might think that a Random Forest Classifier would be a promising candidate for a task such as this.

First, (Treeratpituk and Giles, 2009)

Lastly, there is significant precedent for use of this system

When we used the Random Forest Classifier, we took a look at the most important features for each model. As expected, “Number of characters in target” and “Number of tokens in target” were the number 1 and 2 spots for importance for both languages. Interestingly, the “Rarity score”, which we developed as a basic measure of the use of rare letters in the word was the 3rd most important feature for English and 4th most important feature for Spanish.

It is important to note that we may be getting the wrong impression from these importance scores, since each feature is measured independently. While there is only one “Number of synonyms” feature, there are 312 “Two-letter infix” features and 1297 “Three-letter infix” features. While any given infix feature may be less significant than these more comprehensive, overall features, the sum total of their effect may be greater.

4 Experiments on development set

Does your idea work as expected? Evaluate on the test set the baseline and the improved system, is it still the case? Identify examples in development data which help showcase why the improved system works better.

5 Learning curves

Plot learning curves for the trainable systems you experiment with. Are some systems better than others when less training data is available?

6 Examples of failed predictions

Identify examples where your improved system fails to predict correctly and propose ideas for future work to address them.

7 Conclusions

what have we learnt from your experiments that could inform future work

References

- Pucktada Treeratpituk and C Lee Giles. 2009. Disambiguating authors in academic publications using random forests. In *Proceedings of the 9th ACM/IEEE-CS joint conference on Digital libraries*, pages 39–48. ACM.