

Sentence Generation using Two Models

Final Project in the course DD2380 at KTH

Group 61

K. Hannesson	J. Jóhannsson	E. Ahlsén	J. Andersson
Agust 20	January 12	BIRTHDATE3	February 10
hannesso@kth.se	jokull@kth.se	edvarda@kth.se	jonand8@kth.se



October 13, 2015

Abstract

Bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla

- The following sections are arranged in the order they would appear in a scientific paper. We think that these sections need to be there and written. However, these are only guidelines and if you think that some of these sections or subsections are irrelevant to you, please feel free to remove them. Similarly, if you want to include more sections or subsections please go ahead. Also feel free to rearrange them according to your convenience, but keeping some common sense (eg. Introduction cannot come after Conclusions).
- *Introduction, Related Works, Experimental Results, Discussions, Summary* are sections that MUST be contained.
- In the section of your *Method*: please do not list your project as log book entries, please talk about the final method you want to present to us. Talk about the method scientifically or technically and not as "I did this..." "Then I tried this..." "this happened...." etc.
- Do not paste any code unless it is very relevant!
- The section *Contributions* is a place to express any difference in contributions. The default assumption is that you all agree that all of you had an equal part to play in the project.
- We suggest that you try to write this as scientifically as possible and not simply like a project report. Good Luck!
- Please remove **this** NOTE section in your final report.

[illegible][illegible]

1.2 Outline

Bla bla bla bla bla bla bla bla Section 2 bla bla bla bla bla bla bla bla
Section 3 bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
Section 4 bla bla bla bla bla bla bla bla Section 5 bla bla bla bla bla

2 Related work

Bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla [?] bla bla bla bla bla bla

3 Our method

The group came up with the idea to compare two different approaches to generating text from a corpus, both including grammar but in different ways. To be able to compare them at the same level both approaches used the Brown corpus and trigrams, with fallback to bigrams allowed. A few smoothing techniques were picked for use in both approaches.

The first approach was to create a model that included both word and grammar information. This was achieved by generating trigrams where each word was a tuple of the word and its associated Part-Of-Speech (POS) tag. A trigram from “the man walked” would be (the, DET), (man, NOUN), (walked, VERB))

The second approach separated grammar and words into two models which were used in sequence to generate text. The grammar model provided what POS tag should be next, and the word model provided a word for the given tags, both using trigrams. This would take “the man walked” and create the trigrams (DET, NOUN, VERB) for the grammar model, and (DET, NOUN, man) and (NOUN, VERB, walked) for the grammar-word model.

The smoothing techniques chosen were:

- Maximum Likelihood Estimate
- Laplace smoothing
- Expected Likelihood Estimate
- Simplified Good-Turing Frequency Estimation

We split into pairs with each pair implementing their model. We decided to go with Python 3 for the implementation because we knew the NLTK package would provide us with the necessary building blocks to construct and test the two models. These included

- Brown corpus
- Treebank Part of Speech Tagger (Maximum entropy)
- Punkt Tokenizer Models
- Mappings to the Universal Part-Of-Speech Tagset
- Conditional Frequency Distribution and Conditional Probability Distribution classes
- classes for each of the above mentioned smoothing techniques

Bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla
bla bla bla bla bla bla bla bla bla bla bla

Each model was used to generate 5 sentences per smoothing method for a total of 40 sentences. 5 sentences were also picked from the Brown corpus. All the sentences were mixed together randomly and then added to a survey which allowed participants to rate each sentence on a scale of 1 to 5 whether it was created by a human or a computer, with 1 representing definitely a human and 5 definitely a computer. Apart from the sentence scoring we only asked if the participants were native English speakers or not. The survey solution chosen was QuestionPro.

[illegible]

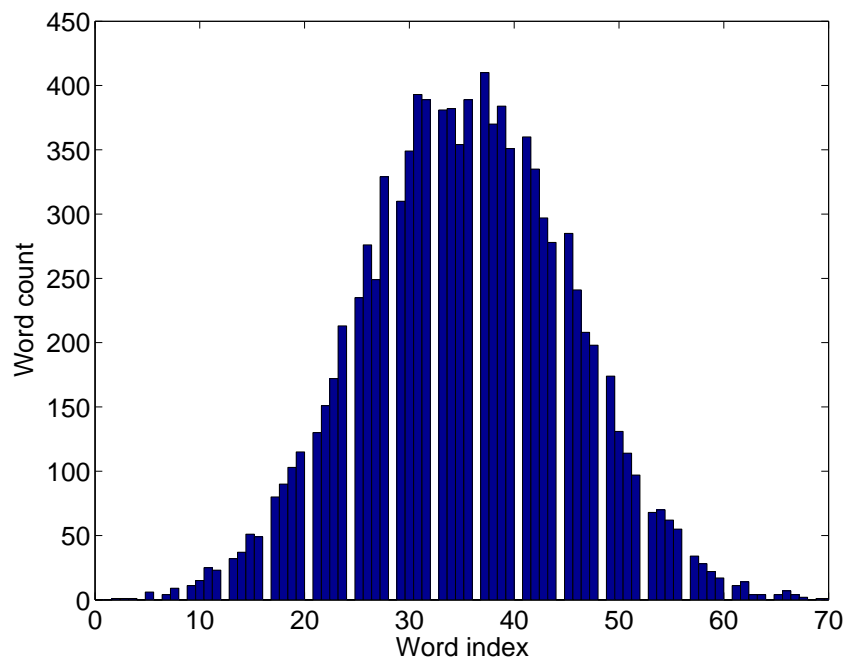


Figure 1: A description that makes browsing the paper easy and clearly describes what is in the picture. Make sure that the text in the figure is large enough to read and that the axes are labelled.

Bla bla	Bla bla	Bla bla
42	42	42
42	42	42

[illegible][illegible]

5 Summary and Conclusions

Bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla bla
bla
bla bla bla bla bla bla bla bla bla bla bla

6 Contributions

We the members of project group 61 unanimously declare that we have all equally contributed toward the completion of this project.