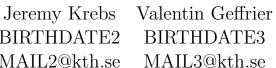
Word prediction performance of n-gram models applied to essentially different corpora

GROUP 34

Sofia Broomé 901210 sbroome@kth.se





Erik Fredriksen BIRTHDATE4 MAIL4@kth.se









Abstract

An investigation of how we can improve word prediction with ngram models applied to both fiction and scientific corpora by adding smoothing techniques and grammar constraints.

NOTE

- 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.

1 Introduction (1–2 pages)

Being able to dissect, classify, analyze and reproduce language is a highly relevant task for various fields. In the realm of artificial intelligence, we want to give language to our agents by means of communicating with them. When we deal with natural language processing we say that we make language models. Seen as there is no finite set of rules that can describe, say, the entire English language in a complete sense, for pragmatic reasons our best option seems to be basing our models on probabilistic observations - regardless of Noam Chomsky's contempt[?] for the notion of probability of a sentence.

At the foundation of every language model that wants to predict words is the concept of n-grams, a method based on probabilistic distributions over length n combinations of subsequent words. An n-gram is a Markov chain of degree n-1. This quite simple construct can capture many patterns in sentences. Even though it doesn't consider grammar explicitly, grammar will inevitably be built in. For instance, an adjective will in many cases be followed by a noun, or a pronoun by a verb, and thus a bigram composed of those two grammatical types in the mentioned order will score high in probability.

An n-gram gives us context for words, albeit not the full one. Gao and Suzuki[?] explore long distance dependency for words through word clusters and the linguistically motivated function word skipping method where function words such as "has", "a", "in", "and", "the", etc, are skipped in favor of more significant words, called head words. In our experiments however, we will not delve further into this subject.

N-grams can also be used in a meta-sense - for instance it's common for part-of-speech-taggers to use n-gram models where they tag the current word based on the last word's tag.

There are some practical issues with the classical n-gram model. What do we do with the n-grams that aren't in our training set and thus have zero probability assigned? This is where techniques of so called smoothing comes in so that our model doesn't fail on encountering a previously unseen word in the test set. In case we are dealing with a higher-order n-gram and we find it has no probability mass, we might want to "back off" from the higher order and estimate the probability for a conditioned unigram, meaning we temporarily look at a smaller portion of a word's history.

Furthermore, what kinds of test sets does our training set allow us to perform well on? One should train on a corpus which is representative of the domain of the intended use. And what happens to our model when we apply it to languages with a higher degree of inflection like Swedish, Basque or German?

From the above examples we see that in many cases just using the n-gram model in itself will not suffice. Over the years researchers in natural language processing have added a lot of tweaks to the original idea such as linear combinations of n-gams, cache language models, LSA-based language models and maximum entropy models, to name a few.

In what follows we will explore n-gram models of varying degrees on dito corpora and grammar to see which results are obtained under which circumstances.

1.1 Contribution

1.2 Outline

2 Related work

3 My method

3.1 Implementation

4 Experimental results

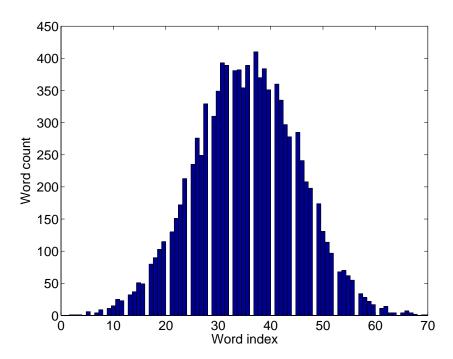


Figure 1: A description that makes browsing the paper easy and clearly describes what is in the picture. Make sure that the text in the figure is large enough to read and that the axes are labelled.

4.1 Experiemntal setup

4.2 Experiment ...

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

Table 1: A description that makes browsing the paper easy and clearly describes what is in the table.

5 Summary and Conclusions

6 Contributions

We the members of project groupXX unanimously declare that we have all equally contributed toward the completion of this project. (PLEASE CHANGE THIS SUITABLY WITH DETAILS, IF IT IS NOT TRUE)