

---

# Billion Word Imputation – Project Proposal

---

## Abstract

In this proposal, we will detail our group's plan to solve the billion word imputation problem. Essentially, given a large training set of sample sentences, we wish to create an algorithm that can predict the missing word in a sentence with one word removed.

## 1. Size of Team

4 Members

## 2. Motivation

There are multiple strong motivations for attempting to solve this task centered around adding a mechanism to accurately recover or replace otherwise undeterminable data. In fields related to auto-correction, progress on this task is useful for correcting writing mistakes on both a word and sentence level. In fields related to auto-suggestion, progress on this task would affect scanning-to-text implementations, damaged textual data set recovery, speech recognition, assisting idiom or phrasing completion, and many other similar issues.

More generally, progress on this task may give insights to language structure by discovering patterns for when phrases or sets of words are strongly associable, and where one has to draw the line between using language knowledge to fill in the blank as opposed to contextual knowledge. Or it's possible that a machine learning based approach has the ability to surpass humans in word imputation, since word imputation in the above contexts can be a difficult enough problem already for humans ourselves.

## 3. Statement of Task

Given a sentence with one word missing, we want to be able to guess what word was removed, and where in the sentence it lies. There are many ways for us to evaluate how good the imputation really is. How close is the location of the imputed word to the location of the missing word? Did we get the correct type of word (noun, adjective, adverb, etc.)? Does our imputed word follow the rules of

basic English grammar? When we propose an imputation, how confident are we that the imputation is correct?

## 4. General Approach

Our first avenue of approach will be to build a Markov Chain model of the training sentences that we would then use to experimentally insert words into the sentence. For every possible location, we will try inserting the words that the Markov Chain model gives a high probability for occurring in that location. If inserting a word in a particular position results in a new sentence that is significantly more likely to be created by the representative Markov Chain than the original sentence was, we may have found the correct insertion. The interesting problems with this approach are to determine the best states to use to build the Markov Chain model. For example, the state of a Markov Chain at a given position can be determined by any number of the preceding words - different length  $n$ -grams ( $n$ -word phrases) can be used as the states of the Markov Chain, and indeed 2-gram or 3-gram states would likely represent the data much better. Combinations of different values of  $n$  can be explored to - is it beneficial to sometimes use 2-grams to determine the probability of the next word, and sometimes use 3-grams?

The main problem of the Markov Chain modeling approach is that language has long-distance dependencies, and the obvious Markov Chain constructions (where states are the previous words) do not. A possible variation on this approach could be keeping track of different words for different lengths of time - when we build the Markov Chain model, the state would consist of whether a word has occurred recently, with 'recently' being defined differently for different words. This idea has aspects of a positional language model, which specifies the probability that two words occur a certain distance from each other in a body of text. Exploration could be done into whether remembering certain types of words (for example, nouns vs. verbs, or subject vs object of the sentence) for longer than other types of words improves the model. However, this introduces further problems. Part-of-speech identification can also be done by Markov Chains<sup>1</sup>, but grammatical relations likely require natural language processing techniques.

---

<sup>1</sup> [http://en.wikipedia.org/wiki/Part-of-speech\\_tagging](http://en.wikipedia.org/wiki/Part-of-speech_tagging)

### 5. Resources

We will be using the test and training data from the Kaggle competition page<sup>2</sup>.

We will also be using parsing libraries for Python or a similar language to parse sentences and grammar. Depending on the computational needs of the algorithm, we may also need to utilize AWS server time to complete the learning in a reasonable amount of time.

### 6. Schedule

*Oct 2 1– Oct 28:* Research background information and methodology for completing the task (Markov chains, n-grams, NLP, etc.)

*Oct 29 – Nov 8:* Create an initial imputation and test it out

*Nov 9 – Nov 11:* Progress report, continue working on imputation

*Nov 15:* Finish some imputation algorithm by this point

*Nov 16 – Nov 26:* Refine existing algorithm to lower test error

*Nov 27 – Dec 4:* Analyze results and complete final poster

---

<sup>2</sup> <https://www.kaggle.com/c/billion-word-imputation/data>