

# Modified Kneser-Ney Language Modeling

Brian Romanowski  
btr@msu.edu

February 26, 2011

# 1 Introduction

Statistical language models are a convenient way to capture linguistic information for use in automated systems. They assign a probability to a sequence of words such that common/expected word sequences have higher associated probability than uncommon/unexpected sequences. Estimating<sup>1</sup> these probabilities from training data allows for good performance in new domains with minimum human intervention. Language models are often combined with other sources of information to improve performance in applications like automated speech recognition, machine translation, and spell checking.

This report describes one particular language modeling technique in detail: the simpler version of Chen and Goodman’s modified Kneser-Ney  $n$ -gram language model [1]. Supporting sections describe the mathematical notation, language modeling background, and the evaluation of language models.

## 2 Notation

Let  $w_t$  indicate the word at index  $t$  in some sequence of  $n$  words  $w_1^n$  chosen from a vocabulary  $V$ . The subsequence starting with the word  $w_i$  and ending with the word  $w_j$ , inclusive, is denoted by  $w_i^j$ . Alternatively, a subsequence of  $n$  words that ends at word  $w_j$  is denoted by  $w_{i-(n-1)}^j$ . This notation is illustrated in figure 1.

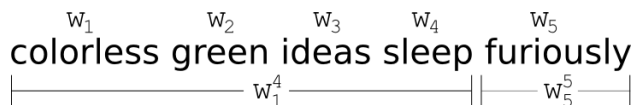


Figure 1: Indexing notation

The probability of a sequence of  $n$  words  $v', v'', \dots, v''^{\dots'}$  is denoted by:

$$P\left(w_{j-(n-1)}^j=v', v'', \dots, v''^{\dots' }\right)=P\left(w_{j-(n-1)}=v', w_{j-(n-2)}=v'', \dots, w_j=v''^{\dots' }\right) \quad (1)$$

In this expression, each  $w_t$  is a categorically distributed random variable with possible outcomes being the finite set of vocabulary words  $v \in V$ .

When the actual word sequence is clear from context or when we make statements about the probability of a generalized word sequence, this simpler notation is used:

$$P\left(w_{j-(n-1)}^j\right)=P\left(w_{j-(n-1)}, w_{j-(n-2)}, \ldots, w_j\right) \quad (2)$$

<sup>1</sup>I would say statistical language models  $\subset$  probabilistic language models; they both report probabilities, but statistical language model parameters are necessarily based on statistical estimates from some training data

The conditional probability of the word at position  $w_j$  given the previous  $n - 1$  word history is used extensively in language modeling:

$$P(w_j \mid w_{j-(n-1)}^{j-1}) = \frac{P(w_{j-(n-1)}^j)}{P(w_{j-(n-1)}^{j-1})} \quad (3)$$

When the conditional probability depends upon words at index  $j < 1$ , such as at the start of a document, generic “words” that serve to mark the start of the document are assumed/faked/inserted.

Counts of word sequences are often used to produce probability estimates like those above. The counting function  $\mathfrak{C}(\cdot)$  counts the number of times word sequences occur, such that:

$$\mathfrak{C}(w_i^j) = d \quad (4)$$

is defined to mean that the word sequence  $w_i^j$  occurred  $d$  times in some sample of text.

The Kneser-Ney family of language model smoothing techniques make use of an additional quantity:

$$N_i(w_{j-(n-1)}^{j-1} \bullet) = \# \text{ unique words that follow } w_{j-(n-1)}^{j-1} \text{ exactly } i \text{ times} \quad (5)$$

The notation  $N_{i+}$  indicates the number of unique words that follow the word sequence at least  $i$  times.

### 3 N-gram Language Models

A straightforward way to estimate the probability of a word sequence is to use the maximum likelihood estimator<sup>2</sup>:

$$\hat{P}(w_i^j) = \frac{\mathfrak{C}(w_i^j)}{N} \quad (6)$$

If the word sequence  $w_i^j$  is relatively long (say six words) compared to the amount of training data that is available (say one billion words of web data), we will not find many occurrences of  $w_i^j$  in the training data. For sequences that occur infrequently, the maximum likelihood estimate  $\hat{P}(w_i^j)$  is not reliable, and will often be a poor estimate of the true probability  $P(w_i^j)$ .

One step towards a solution is to make the word sequences short, so that they are more likely to occur often in training data. To accomplish this process, we first decompose the probability of a sequence of words according to the chain rule of probability:

$$\begin{aligned} P(w_i^j) &= P(w_i^{j-1}) \cdot P(w_j \mid w_i^{j-1}) \\ &= P(w_i) \cdot P(w_{i+1} \mid w_i) \cdots P(w_{j-1} \mid w_i^{j-2}) \cdot P(w_j \mid w_i^{j-1}) \end{aligned} \quad (7)$$

If we assume that each word  $w_j$  is only dependent on the previous  $n - 1$  words, equation 7 simplifies to the following:

$$P(w_i^j) \approx P(w_j \mid w_{j-(n-1)}^{j-1}) \cdot P(w_{j-1} \mid w_{j-(n-1)}^{j-2}) \cdots P(w_{i+1} \mid w_i) \cdot P(w_i) \quad (8)$$

<sup>2</sup>This is the maximum likelihood estimator for a multinomial distribution over conditionally independent categorical events. It is analogous to the maximum likelihood estimator for a binomial distribution over a sequence of Bernoulli random variables.

This is the *n*-gram assumption.

In an *n*-gram model, the probability of the  $n^{\text{th}}$  word depends only upon the  $n - 1$  previous words. Each of these conditional probability distributions (i.e., model parameters) may be estimated by the relative frequency calculation:

$$\hat{P}(w_j | w_{j-(n-1)}^{j-1}) = \frac{\hat{P}(w_{j-(n-1)}^j)}{\hat{P}(w_{j-(n-1)}^{j-1})} \quad (9)$$

$$= \frac{\mathfrak{C}(w_{j-(n-1)}^j)/N}{\mathfrak{C}(w_{j-(n-1)}^{j-1})/N} \quad (10)$$

$$= \frac{\mathfrak{C}(w_{j-(n-1)}^j)}{\mathfrak{C}(w_{j-(n-1)}^{j-1})} \quad (11)$$

*N*-gram models are often referred to with Latin prefixes for low order *n*-grams and with numeral prefixes for high order *n*-grams as: unigram, bigram, trigram, 4-gram, 5-gram, etc.

## 4 The Problem of Sparsity

The *n*-gram assumption helps mitigate sparsity problems, but there are still many word sequences that are seen rarely or not at all in training data. The worst case occurs when there is a sequence of words encountered during testing that was not encountered during training.

As an example, consider a bigram model is trained on the Brown corpus. Suppose the model is used to calculate the probability of the sentence: “the quick brown fox”.

$$P(w_1^4) = P(\text{“the”})P(\text{“quick”}|\text{“the”})P(\text{“brown”}|\text{“quick”})P(\text{“fox”}|\text{“brown”}) \quad (12)$$

It is reasonable that all of these bigrams will appear in a large training corpus and will have non-zero probabilities. However, it may be the case that some bigrams in the similar sentence “the quick *wily* fox” will *not* appear in training data.<sup>3</sup>

If the bigram “quick wily” was not seen in the training data, its relative frequency estimate is zero.

$$P(w_1^4) = P(\text{“the”})P(\text{“quick”}|\text{“the”})P(\text{“wily”}|\text{“quick”})P(\text{“fox”}|\text{“wily”}) \quad (13)$$

$$= P(\text{“the”})P(\text{“quick”}|\text{“the”}) \cdot 0 \cdot P(\text{“fox”}|\text{“wily”}) \quad (14)$$

$$= 0 \quad (15)$$

The presence of just one unseen bigram can prevent consideration of a syntactically and semantically valid sentence.

## 5 Modified Kneser-Ney Smoothing

Kneser-Ney smoothing [3] was modified slightly by Chen and Goodman [1]. Below is a description of the simpler version of their algorithm; the more complicated version that optimizes parameter values on held-out data is not discussed.

<sup>3</sup>All bigrams in 12 appear at least 100,000 times in reported Google search results, whereas the bigram “quick wily” appears in only 500 reported Google search results.

The modified Kneser-Ney smoothing relation is defined recursively as follows:

$$\hat{P}(w_j \mid w_{j-(n-1)}^{j-1}) = \frac{\mathfrak{C}(w_{j-(n-1)}^j) - D_n(\mathfrak{C}(w_{j-(n-1)}^j))}{\sum_{w_j} \mathfrak{C}(w_{j-(n-1)}^{j-1}, w_j)} + \gamma(w_{j-(n-1)}^{j-1}) \hat{P}(w_j \mid w_{j-(n-2)}^{j-1}) \quad (16)$$

where

$$D_n(c) = \begin{cases} 0 & \text{if } c = 0 \\ 1 - 2 \left( \frac{n_1}{n_1 + 2n_2} \right) \left( \frac{n_2}{n_1} \right) & \text{if } c = 1 \\ 2 - 3 \left( \frac{n_1}{n_1 + 2n_2} \right) \left( \frac{n_3}{n_2} \right) & \text{if } c = 2 \\ 3 - 4 \left( \frac{n_1}{n_1 + 2n_2} \right) \left( \frac{n_4}{n_3} \right) & \text{if } c \geq 3 \text{ indicated by “} c = 3+ \text{”} \end{cases} \quad (17)$$

where  $n_i$  are the number of  $n$ -grams of order  $n$  that appear  $i$  times in the training data, and:

$$\gamma(w_{j-(n-1)}^{j-1}) = \frac{D_n(1)N_1(w_{j-(n-1)}^{j-1} \bullet) + D_n(2)N_2(w_{j-(n-1)}^{j-1} \bullet) + D_n(3+)N_{3+}(w_{j-(n-1)}^{j-1} \bullet)}{\sum_{w_j} \mathfrak{C}(w_{j-(n-1)}^{j-1}, w_j)} \quad (18)$$

We assume a closed vocabulary and terminate the recursion such that the unigram probability is just the maximum likelihood estimate:

$$\hat{P}(w_j) = \frac{\mathfrak{C}(w_j)}{N} \quad (19)$$

These equations are all taken directly from the Chen and Goodman paper; however, we make explicit the fact that the parameters  $D_n(c)$  are calculated for each order  $n$ .

Things to prove or examine (most of these things are proven in Chen and Goodman):

- $\sum_{w_j} \hat{P}(w_j \mid w_{j-(n-1)}^{j-1}) = 1$
- intuition of  $\gamma(\cdot)$
- intuition behind KN
- Is this generally how SRILM does it?

## 6 Evaluation

For a particular task, language model performance is best modeled by using a language model as a component *in situ*, in a full system. In this setup, the language model may effect system accuracy, but may also be analyzed with respect to training time, model size, and performance as the task domain changes in time. However, the necessity of repeated system runs and the difficulty of collecting adequate test data can make *in situ* evaluations costly. Additionally, it is not obvious how performance results conditioned on a particular task informs the *in situ* performance in a different task or domain.

## 6.1 Perplexity

The *perplexity* metric measures language model performance with respect to textual testing data. Language model perplexity is often strongly correlated with task performance. Perplexity may be defined as shown [2]:

$$\text{PP}(W) = P(w_1^N)^{-\frac{1}{N}} \quad (20)$$

where  $\text{PP}(W)$  is the perplexity of a test data set  $w_1^N$ .

This definition of perplexity may be interpreted as the geometric mean of the inverse probability of the test set. The probability of the test data is a product of  $N$  individual conditional  $n$ -gram probabilities. The geometric mean is the number that, when multiplied  $N$  times, produces an equivalent probability. This value characterizes the average uncertainty of a word in the testing data.

The perplexity, as the inverse of the geometric mean, can be viewed as average size of the set of equiprobable words at each index in the test dataset. Thus, *lower perplexity indicates a better fitting model*, because there is a smaller number of word alternatives for each word in the testing data, on average.

Perplexity is also related to the information theoretic quantity *cross entropy*. See appendix 7 for the proof that perplexity is related to cross entropy in the following manner:

$$H(p, q) = \lg \text{PP}(w_1^N) \quad (21)$$

where  $H(\cdot, \cdot)$  is cross entropy and  $\lg$  is the base 2 logarithm.

Improvements in perplexity do not always correspond to improvements in task performance. Perplexity is evaluated only on valid sentence test data, so no measure of false positive performance is indicated. Additionally, a language model that assigns more probability mass to common words may increase perplexity, but it may also “overwhelm” the probability mass contributed by other sources of information, like the acoustic model in speech recognition decoders.

## 6.2 Statistical Significance

# 7 Perplexity and Cross Entropy

An argument will be made for use of the cross entropy metric when evaluating language model data. The cross entropy metric is also shown to be equivalent to the previously defined perplexity metric. The reasoning presented below is taken from Jurafsky and Martin [2]; additional annotation and discussion is also provided.

First, consider the entropy of a language consisting of all possible word sequences, as the length of the word sequences approaches infinity.

$$H(p) = - \lim_{n \rightarrow \infty} \frac{1}{n} H_n \quad (22)$$

$$= - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{w_1^n} p(w_1^n) \lg p(w_1^n) \quad (23)$$

$$= - \lim_{n \rightarrow \infty} \frac{1}{n} \lg p(w_1^n) \quad (24)$$

where  $\lg x$  is the base-2 logarithm of  $x$ ,  $H(p)$  is the total entropy of a language  $W$ ,  $H_n$  is the entropy of a sequence of words of length  $n$ , and  $w_1^n$  is the set of all sequences of length  $n$ . In equations 24, we make use of the Shannon-McMillan-Breiman theorem. The necessary assumption is that the language  $W$  is ergodic, meaning that all possible word sequences appear in an infinitely long sample from the language, and that

each word sequence appears with a frequency corresponding to its true probability of appearance. Then, looking at an infinitely long string captures the same information available from any weighted, shorter string.

Next, we examine the cross entropy  $H(p, q)$ :

$$H(p, q) = - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{W_1^n} p(W_1^n) \lg q(W_1^n) \quad (25)$$

$$= - \lim_{n \rightarrow \infty} \frac{1}{n} \lg q(W_1^n) \quad (26)$$

$$\approx - \frac{1}{N} \lg q(W_1^N) \quad (27)$$

where the Shannon-McMillan-Breiman theorem is used again in equation 26. The approximation in equation 27 is derived by assuming that the testing data on which we wish to evaluate a language model is long enough to capture important language properties.

This derivation allows us to state that the entropy of a language is bounded by the cross entropy of our language model:

$$H(p) \leq H(p, q) \quad (28)$$

It should be clear that a better language model will approach “closer”, that is have lower cross entropy, to the entropy of the language.

Perplexity is an equally valid metric, as it is related to cross entropy as follows:

$$H(p, q) = \lg PP(W_1^N) \quad (29)$$

## 8 License

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

## References

- [1] S. Chen and J. Goodman. An empirical study of smoothing techniques for language modeling. Technical Report TR-10-98, Harvard, 1998.
- [2] D. Jurafsky and J. Martin. *Speech and Language Processing*. 2 edition, 2008.
- [3] R. Kneser and H. Ney. Improved backing-off for m-gram language modeling. In *Acoustics, Speech, and Signal Processing, 1995. ICASSP-95., 1995 International Conference on*, volume 1, pages 181–184, 1995.