

On the Use of Morphological Constraints in N-gram Statistical Language Model

A. Ghaoui⁽¹⁾⁽²⁾⁽³⁾, F. Yvon⁽²⁾, C. Mokbel⁽¹⁾ et G. Chollet⁽²⁾

⁽¹⁾ University of Balamand BP 100 Tripoli, Lebanon
Chafic.Mokbel@balamand.edu.lb

⁽²⁾ CNRS-URA820, Ecole Nationale Supérieure des Télécommunications
46 rue Barrault, 75634 Paris cedex 13, France
{Francois.Yvon, Gerard.Chollet}@enst.fr

⁽³⁾ Jinny Software Ltd, Samra Center, 1st floor, Fanar, El Metn, Lebanon
Antoine.Ghaoui@jinny.ie

Abstract

State of the art Speech Recognition systems use statistical language modeling and in particular N-gram models to represent the language structure. The Arabic language has a rich morphology, which motivates the introduction of morphological constraints in the language model. Class-based N-gram models have shown satisfactory results, especially for language model adaptation and training from reduced datasets. They were also proven quite effective in their use of memory space. In this paper, we investigate a new morphological class-based language model. Morphological rules are used to derive the different words in a class from their stem. As morphological analyzer, a rule-based stemming method is proposed for the Arabic language. The language model has been evaluated on a database composed of articles from Lebanese newspaper Al-Nahar for the years 1998 and 1999. In addition, a linear interpolation between the N-gram model and the morphological model is also evaluated. Preliminary experiments detailed in this paper show satisfactory results.

1. Introduction

Language modeling aims at capturing local syntactic constraints between words. Statistical language models represent these constraints as probability distributions over words sequences. In the case of N-gram language models, the probability of a word sequence is computed as a product of the conditional probabilities of each word given a restricted history or context as follows:

$$P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{i-1}, \dots, w_{i-N+1}) \quad (1)$$

For the i^{th} word, the context/history is limited to the previous N-1 words. Even with such restricted contexts, N-gram models include a huge number of parameters. A bi-gram (N=2) model on a vocabulary of V words has V^2 parameters $P(w_i|w_j)$, where $P(w_i|w_j)$ represents the probability to have word j followed by word i . The estimation and use of such a large parameter set raise several practical problems. First, the training corpus has to be large enough to estimate all these parameters. In fact, even if a very large corpus were available, not all the combinations would appear, since some of them are extremely rare. Secondly, N-gram models require a large memory space.

In order to overcome these problems, a parametric distribution function with a limited number of parameters is generally used

for the conditional probability distribution of the vocabulary words. However, the use of a parametric distributions function requires a distance measure or at least an order to be defined on the set of vocabulary words. No such distance measure or order exists. A possible solution is to cluster the vocabulary words into classes and to consider that words in a given class share certain properties. The use of word classes smoothes the probability distribution of the (word based) N-gram model. Syntactical classes, semantic classes, morphological classes or classes induced directly from the data have been considered. In this work, we are particularly interested in morphological constraints as a basis for class-based N-gram modeling. In our model, we consider a word to be uniquely derived from its root by application of a morphological rule. The probability of a word, given its context, thus combines two terms: one is the contribution of the root based N-gram model; the other is probability of application of the morphological rule. This model is introduced in Section 2.

Section 3 describes the morphological model where a transducer is used to generate the different rules. The rules are empirical. This morphological class-based model has been linearly combined to the classical N-gram model.

Our morphological class-based language model has been implemented within the SRILM toolkit [1]. Preliminary experiments have been conducted using the articles of newspaper Al-Nahar for the years 1998 and 1999 as database. The database, experiments and results are detailed in Section 4. The final section of the paper presents conclusions and perspectives.

2. Morphological class based Language Model

Class-based language models have been studied (see e.g. [2; 3; 4]). In this paper, a particular approach is proposed where morphological rules are used to define the classes. Morphologically rich languages have a large set of words that derive from a small set of roots. The basic idea in our morphological class based model is to define a class per root and to associate all derived words to that class. The underlying assumption is that a strong relationship exists between the roots of the words in a language. The meaning of a sentence is completed by adding to the roots the derivation rules information.

Let w_i be a word of the vocabulary, r_i the root of w_i and g_i the morphological rule which allows the word w_i to be derived from r_i . At this stage, we assume that a word derives from a unique root, which may not be true in some cases. Each possible root in the vocabulary defines a class, which contains all the words deriving from that root.

These results show that the morphological class-based model yields poor performance compared to a classical trigram. This is in line with results generally obtained with class-based N-gram models, and is due to the fact that the class-based model has a much smaller number of parameters. These results are expected, based on the different simplifications applied in our model of (Eq. 6).

4.4. Linear Interpolation

In order to investigate if the morphological class-based model may be combined to a classical N-gram model, linear interpolation has been used. The linear interpolation permits a reduction in the perplexity by a factor of 1.7% (369.893) as shown in the figure 2. This shows that our morphological class-based model can bring some improvement over the classical N-gram model.

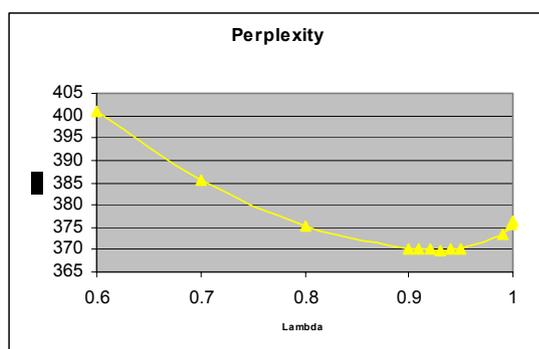


Figure 2: Results with linear interpolation between the morphological class-based model and the classical trigram model.

5. Conclusions and Perspective

In this paper an N-gram statistical language model that integrates morphological constraints in the form of classes is proposed. A complete theoretical framework is provided to integrate morphological constraints in the N-gram modeling. Depending on the assumptions made, multiple models may be derived from our theoretical framework.

A simplified model based on morphological classes has been applied to the Arabic language. A simple Arabic morphological analyzer has been developed for this purpose where empirical morphological rules are used in the form of a transducer. The SRILM toolkit has been modified to integrate the simplified model. Preliminary tests show a degradation of the performance in term of perplexity when compared to the classical N-gram. This was expected due to the reduction in the number of parameters in the modeling. The loss in term of perplexity is compensated by a reduction in the complexity of the model. The linear interpolation between the classical N-gram and the model based on morphological class improves the result in term of perplexity. This has been shown in our experiments.

This work will be pursued in various directions. First a more complete and accurate morphological analyzer will be used. Second, the morphological class-based language model will be experimented with fewer assumptions. Moreover, the assumption that a unique stem exists for a word will be

relaxed by considering several stems per word with a probability associated to each of them.

6. Acknowledgements

This work is partially supported by the EC project NEMLAR.

7. References

- [1] Stolcke A., "SRILM - An Extensible Language Modeling Toolkit," *Proc. Intl. Conf. Spoken Language Processing*, Denver, Colorado 2002.
- [2] Manning C. and Schutze H, *Foundations of Statistical Natural Language Processing*, MIT Press. Cambridge 1999.
- [3] Brown P. F., Della Pietra V. J., deSouza P. V., Lai J. C., and Mercer R. L., Class-based N-gram models of natural language, *Computational Linguistics*, Vol. 18, pp. 467-479 1992
- [4] Neisler T, Category-based statistical language models, *Ph.D thesis, Department of Engineering, University of Cambridge, U.K., 1997*
- [5] Darwish K, "Building a Shallow Arabic Morphological Analyzer in One Day" *Proceedings of the ACL Workshop on Computational Approaches to Semitic Languages*, Philadelphia, Boston, 2002.
- [6] Beesley K, "Arabic Finite-State Morphological Analysis and Generation," *COLING-96*, Vol. 1, pp. 89-94 1996
- [7] Katz S, "Estimation of Probabilities from Sparse Data for the Language Model Component of a speech Recognizer", *IEEE Trans. on Acoustics, Speech, and Signal Processing*, Vol. ASSP-35, no. 3, pp. 400-401 MARCH 1987