

Bayesian Nets in Syntactic Categorization of Novel Words

Leonid Peshkin

Dept. of Computer Science
Harvard University
Cambridge, MA
pesh@eecs.harvard.edu

Avi Pfeffer

Dept. of Computer Science
Harvard University
Cambridge, MA
avi@eecs.harvard.edu

Virginia Savova

Dept. of Cognitive Science
Johns Hopkins University
Baltimore, MD
savova@jhu.edu

Abstract

This paper presents an application of a Dynamic Bayesian Network (DBN) to the task of assigning Part-of-Speech (PoS) tags to novel text. This task is particularly challenging for non-standard corpora, such as Internet lingo, where a large proportion of words are unknown. Previous work reveals that PoS tags depend on a variety of morphological and contextual features. Representing these dependencies in a DBN results in an elegant and efficient PoS tagger.

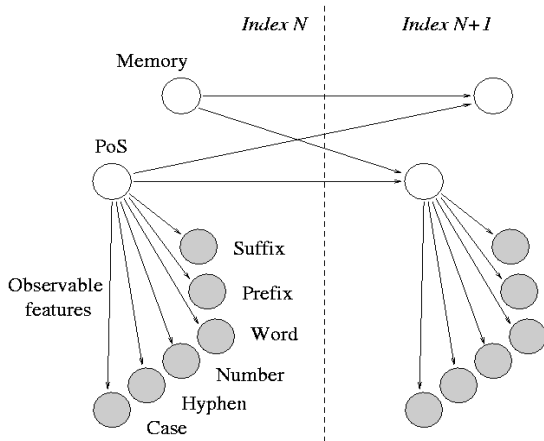
1 Introduction

Uncovering the syntactic structure of a text is a necessary step towards extracting its meaning. In order to obtain an accurate parse for an unseen text, we need to assign Part-of-Speech (PoS) tags to a string of words. This paper covers one aspect of our work on Dynamic Bayesian Networks (DBNs) for PoS tagging, which demonstrates their success at tagging out-of-vocabulary (OoV) words. Please refer to the companion paper [Peshkin'03] for substantial discussion of our method and other details. Although currently existing algorithms exhibit high word-level accuracy, PoS tagging is not a solved problem. First, even a small percentage of errors may derail subsequent processing steps. Second, the results of tagging are not robust if a large proportion of words are unknown, or if the testing corpus differs in style from the training corpus. At the same time, diverse training corpora are lacking and most taggers are trained on a large annotated corpus extracted from the Wall Street Journal (WSJ).

These factors significantly hamper the use PoS tagging to extract information from non-standard corpora, such as email messages and websites. Our work on Information Extraction from an email corpus [Peshkin&Pfeffer'03] left us searching for a PoS tagger that would perform well on Internet texts and integrate easily into a large probabilistic reasoning system by producing a distribution over tags rather than deterministic answer.

Internet sources exhibit a set of idiosyncratic characteristics not present in the training corpora available to taggers to date. They are often written in telegraphic style, omitting closed-class words, which leads to a higher percentage of ambiguous items. Most importantly, as a consequence of the rapidly evolving Netlingo, Internet texts are full of new words, misspelled words and one-time expressions. These characteristics are bound to lower the accuracy of existing taggers.

A look at the literature confirms that error rates for unknown words are quite high. According to several recent publications [Toutanova&Manning '02, Lafferty et al.'02] OoV tagging still presents a serious challenge to the field. The transformation-based Brill tagger [Brill'95], achieves 96.5% accuracy for the WSJ, but a mere 85% on unknown words. Existing probabilistic taggers also don't fare well on unknown words. Reported results on OoV rarely exceed Brill's performance by a tiny fraction. They are mostly based on (Hidden) Markov Models [Brants'00, Kupiec'92]. A model based on Conditional Random Fields [Lafferty et al.'02] outperforms the HMM tagger on unknown words yielding 24% error rate. The best result known to us [Toutanova&Manning'02] is achieved by enriching the feature representation of the MaxEnt approach [Ratnaparkhi, 1996].



2 Bayesian Net for PoS Tagging

Our model is deliberately based on the original feature set of Ratnaparkhi's MaxEnt (unlike Toutanova&Manning[2002]). A set of binary features and a set of vocabulary features are combined into a Bayesian network. The binary features indicate the presence or absence of a particular character in the token: does the token contain a **capital** letter; does the token contain a **hyphen**; does the token contain a **number**. We used Ratnaparkhi's vocabulary lists to encode the values of 6458 frequent **Words**, 3602 **Prefixes** and 2925 **Suffixes** up to 4 letters long.

A dynamic Bayesian network (DBN) is a Bayesian network unwrapped in time, such that it can represent dependencies between variables at adjacent positions. Murphy [2002] gives a good overview of DBNs. The set of observable variables in our network consists of the binary and vocabulary features mentioned above. In addition, there are two hidden variables: PoS and Memory which reflects contextual information about past PoS tags. Unlike Ratnaparkhi, we do not directly consider any information about preceding words. However, a special value of Memory indicates whether we are at the beginning of the sentence.

Learning in our model is equivalent to collecting statistics over co-occurrences of feature values and tags. It is implemented as GAWK scripts and takes just minutes on the WSJ training corpus. Compare this to a laborious Improved Iterative Scaling for MaxEnt. Tagging is carried out by the standard Forward-Backward algorithm (see e.g. Murphy[2002]). There is no need to use specialized search algorithms such as Ratnaparkhi's "beam search". In addition, our method does not require a "Development" stage.

Following the established data split we use sections (0-22) of WSJ [Marcus'94] for training and the rest (23-24) as a test set. The test sections contain 4792 sentences out of about 55600 total sentences in WSJ corpus, an average of 23 tokens per sentence. In addition, we created two specialized testing corpora (available upon request for comparison purposes). A small Email corpus was made of excerpts from the MUC seminar announcement corpus. "The Jabberwocky" is a poem by Lewis Carroll where about one third of the words are made-up, but their syntactic categories are apparent to speakers of English. We use "The Jabberwocky" to illustrate performance on unknown words. Both the Email corpus and the Jabberwocky were pre-tagged by the Brill tagger and then manually corrected.

We began our experiments by using the original set of features and vocabulary lists of Ratnaparkhi for the variables Word, Prefix and Suffix. This produced a reasonable result. However, while investigating the relative contribution of each feature in this setting, we discovered that the removal of the three binary features from the feature set does not significantly alter performance. Upon close examination, the vocabularies turned out to contain a lot of redundant information that is otherwise handled by these features. For example, Prefix list contained 84 hyphens (e.g. both "co-" and "co"), 530 numbers and 1500 capitalized words, including capital letters. We proceed, using reduced vocabularies obtained by removing redundant information from the original lists. The results for various testing conditions are presented in a table.

Description	Average	OoV	Sentence
Original feature set of Ratnaparkhi	6.8	13.2	69.4
Email corpus	16.3	12.2	79.0
Jabberwocky	11.0	23.0	65.0
Trained on WSJ tested on Brown	13.1	26.5	73.2
Factored feature Set on random WSJ	3.6	9.8	52.7
Factored feature set on WSJ 23-24	3.6	9.4	51.7

Our average error rate is comparable to the best result known on this benchmark (e.g. Toutanova&Manning[2002]). At the same time, our performance on OoV words is significantly better (9.4% versus 13.3%). We attribute this difference to the purer representation of morphologically relevant suffixes in our factored vocabulary, which excludes redundant and therefore potentially confusing information. Another reason may be that our method puts a greater emphasis on the syntactically relevant facts, such as morphology and tag sequence information by *refraining* to use word-specific cues. We also found that removing prefix information completely did not worsen performance (*contra* Toutanova&Manning’02).

Despite our good performance on the WSJ corpus, we failed to improve Brill’s tagging on our two specialized corpora. Both Brill and our method achieved 89% on the Jabberwocky poem. Note, however, that Brill uses much more sophisticated mechanisms to obtain this result. It was particularly disappointing for us to find out that we did not succeed in labeling the Email corpus accurately (16.3% versus 14.9% of Brill). However, the reason for this poor performance appears to be partly related to a labeling convention of the Penn Treebank, which essentially causes most capitalized words to be categorized as NNPs. In our view, there is a significant difference between the grammatical status of such proper names as “Bill Gates”, where words can not be said to modify one another, and a name of an institution such as “Department of Electrical Engineering”, where “electrical” clearly modifies “engineering”. While a rule-based system profits from this simplistic convention, our method is harmed by it.

3 Conclusion

Our approach shows promise as it is both probabilistic and outperforms existing statistical taggers on unknown words. We are especially encouraged by our performance on the WSJ and take this as evidence that our method has the potential to significantly improve PoS tagging of non-standard texts. In addition, our method has the advantage of being conceptually simple, fast, and flexible with respect to feature representation. We

are currently investigating the performance of other DBN topologies on PoS tagging.

References

- T. Brants. 2000. *TnT - a statistical part-of-speech tagger*. In Proceedings of the 6th ANLP.
- E. Brill. 1995. *Transformation-based error-driven learning and natural language processing*. Computational Linguistics, 21(4):543—565.
- E. Charniak, C. Hendrickson, M. Jacobson and M. Perkowski. 1993. *Equations for part-of-speech tagging*. In Proceedings of 11th AAAI.
- F. Jelinek. 1985. *Markov source modeling of text generation*. In J. Skwirzinski, ed., Impact of Processing Techniques on Communication, Dordech.
- J. Kupiec. 1992. *Robust part-of-speech tagging using a hidden Markov model*. Computer Speech and Language, 6:225-242.
- J. Lafferty, A. McCallum and F. Pereira. 2002. *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data*, In Proceedings of 18th ICML.
- M. Marcus, G. Kim, M. Marcinkiewicz, R. MacIntyre, A. Bies, M. Ferguson, K. Katz and B. Schasberger. 1994. *The Penn Treebank: Annotating predicate argument structure*. ARPA Human Language Technology Workshop.
- C. Manning and H. Schutze. 1999. *Foundations of Statistical Natural Language Processing*. The MIT Press. Cambridge, MA.
- K. Murphy. 2002. *Dynamic Bayesian Networks: Representation, Inference and Learning*. PhD thesis. UC Berkeley.
- L. Peshkin and A. Pfeffer. 2003. *Part-of-Speech tagging with Dynamical Bayesian Network*. manuscript.
- L. Peshkin and A. Pfeffer. 2003. *Bayesian Information Extraction Network*. manuscript.
- A. Ratnaparkhi. 1996. *A maximum entropy model for part-of-speech tagging*. In Proceedings of EMNLP.
- C. Samuelsson. 1993. *Morphological Tagging Based Entirely on Bayesian Inference*. 9th Nordic Conference on Computational Linguistics, Stockholm University, Stockholm, Sweden.
- K. Toutanova and C. Manning. 2002. *Enriching the Knowledge Sources Used in a Maximum Entropy PoS Tagger*. Proceedings of the 6th CoNLL.