

Named Entity Recognition: Exploring Features

Maksim Tkachenko

St Petersburg State University
St Petersburg, Russia

maksim.tkachenko@math.spbu.ru andrey.simanovsky@hp.com

Andrey Simanovsky

HP Labs Russia
St Petersburg

Abstract

We study a comprehensive set of features used in supervised named entity recognition. We explore various combinations of features and compare their impact on recognition performance. We build a conditional random field based system that achieves 91.02% F_1 -measure on the CoNLL 2003 (Sang and Meulder, 2003) dataset and 81.4% F_1 -measure on the OntoNotes version 4 (Hovy et al., 2006) CNN dataset, which, to our knowledge, displays the best results in the state of the art for those benchmarks respectively. We demonstrate statistical significance of the boost of performance over the previous top performing system. We also obtained 74.27% F_1 -measure on NLPBA 2004 (Kim et al., 2004) dataset.

1 Introduction

Recognition of named entities (e.g. people, organizations, locations, etc.) is an essential task in many natural language processing applications nowadays. Named entity recognition (NER) is given much attention in the research community and considerable progress has been achieved in many domains, such as newswire (Ratinov and Roth, 2009) or biomedical (Kim et al., 2004) NER. Supervised NER that uses machine learning algorithms such as conditional random fields (CRF) (McCallum and Li, 2003) is especially effective in terms of quality of recognition.

Supervised NER is extremely sensitive to selection of an appropriate feature set. While many features were proposed for use in supervised NER

systems (Krishnan and Manning, 2006; Finkel and Manning, 2009; Lin and Wu, 2009; Turian et al., 2010), only limited studies of the impact of those features and their combinations on the effectiveness of NER were performed. In this paper we provide such a study.

Our contributions are the following:

- analysis of the impact of various features taken from a comprehensive set on the effectiveness of a supervised NER system;
- construction of a CRF-based supervised NER system that achieves 91.02% F_1 -measure on the CoNLL 2003 (Sang and Meulder, 2003) dataset and 81.4% F_1 -measure on the OntoNotes version 4 (Hovy et al., 2006) CNN dataset;
- demonstration of statistical significance of the obtained boost in NER performance on the benchmarks;
- application to NER of a DBpedia (Mendes et al., 2011) markup feature and a phrasal clustering (Lin et al., 2010) feature which have not been considered for NER in previous works.

The remainder of the paper is structured in the following way. In Section 2 we describe related work on feature analysis. In Section 3 we give a brief introduction to the benchmarks that we use. In Section 4 we discuss various features and their impact. Section 5 describes the final proposed system. Section 6 contains a summary of the performed work and future plans.

2 Related Work

The majority of papers on NER describe a particular method or feature evaluation and do not make a systematic comparison of combinations of features. In this paper those works are mentioned later when we discuss a particular feature or a group of features. In this section we present several works that deal with multiple features and thus are close to our study.

Design questions of NER systems were considered by (Ratinov and Roth, 2009). They used a perceptron-based recognizer with greedy inference and evaluated two groups of features: non-local dependencies (e.g. context aggregation) and external information (e.g. gazetteers mined from Wikipedia). Their recognizer was tested on the CoNLL 2003 dataset, a newswire dataset ($F_1 = 90.80\%$), the MUC7 dataset, and their own web pages dataset.

The authors of (Turian et al., 2010) systematically compared word representations in NER (Brown clusters, Collobert and Weston embeddings, HLBL embeddings). They ignored other types of features.

(Saha et al., 2009) presented a comparative study of different features in biomedical NER. They used a dimensionality reduction approach to select the most informative words and suffixes and they used clustering to compensate for the lost information. The MaxEnt tagger developed by them obtained $F_1 = 67.4\%$ on NLPBA 2004 data.

3 Benchmarks

In this paper we present the results obtained on three benchmarks: CoNLL 2003, OntoNotes version 4, and NLPBA 2004 dataset.

CoNLL 2003 is an English language dataset for NER. The data comprises Reuters newswire articles annotated with four entity types: person (PER), location (LOC), organization (ORG), and miscellaneous (MISC). The data is split into a training set, a development set (testa), and a test set (testb). Performance on this task is evaluated by measuring precision and recall of annotated entities combined into F_1 -measure. We used BILOU (begin, inside, last, outside, unit) annotation scheme to encode named entities. Previous top performing systems also followed that scheme. We study feature behavior on this benchmark; our system is tuned on the test and development sets of it.

OntoNotes version 4 is an English language dataset designed for various natural language processing tasks including NER. The dataset consists of several sub-datasets taken from different sources including Wall Street Journal, CNN news, machine-translated Chinese magazines, Web blogs, etc. We provide the results obtained by our final system on OntoNotes subsets in order to compare them with earlier works. It has its own set of named entity classes but it has a mapping of those to CoNLL classes. We use the latter for systems comparison. We used the same test/training split as in (Finkel and Manning, 2009).

NLPBA 2004 dataset (Kim et al., 2004) is an English language dataset for bio-medical NER. It consists of a set of PubMed abstracts nad has a corresponding set of named entites (protein, DNA, RNA, cell line, cell type).

4 Feature Set

We performed feature comparison using our system which is a CRF with Viterbi inference. We have also tested greedy inference and have found out that the system performs worse and its results are lower than those of a perceptron with greedy inference that we modeled after (Ratinov and Roth, 2009).

In each of the following subsections we consider a particular type of features. In Subsection 4.1 we deal with local knowledge features which can be extracted from a token (word) being labeled and its surrounding context. Subsection 4.2 describes evaluation of external knowledge features (part-of-speech tags, gazetteers, etc.). Discussion of non-local dependencies of named entities is included in Subsection 4.3. Subsection 4.4 contains further improvements of performance and specific features that do not fall into previous categories; they help to overcome common errors on the CoNLL 2003 dataset.

4.1 Local Knowledge

Our baseline recognizer uses only local information about a current token. It is not surprising that a token-based CRF performs poorly, espe-

References

- Rie Kubota Ando and Tong Zhang. 2005. A high-performance semi-supervised learning method for text chunking. In Kevin Knight, Hwee Tou Ng, and Kemal Oflazer, editors, *ACL*. The Association for Computer Linguistics.
- Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören Auer, Christian Becker, Richard Cyganiak, and Sebastian Hellmann. 2009. Dbpedia - a crystallization point for the web of data. *J. Web Sem.*, 7(3):154–165.
- Peter F. Brown, Vincent J. Della Pietra, Peter V. Desouza, Jennifer C. Lai, and Robert L. Mercer. 1992. Class-Based n-gram Models of Natural Language. *Computational Linguistics*, 18(4):467–479.
- Christos Christodoulopoulos, Sharon Goldwater, and Mark Steedman. 2010. Two decades of unsupervised pos induction: how far have we come? In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, EMNLP ’10, pages 575–584, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Grzegorz Chrupala. 2011. Efficient induction of probabilistic word classes with LDA. In *Proceedings of 5th International Joint Conference on Natural Language Processing*, pages 363–372, Chiang Mai, Thailand, November. Asian Federation of Natural Language Processing.
- M. Ciaramita and Y. Altun. 2005. Named-entity recognition in novel domains with external lexical knowledge. In *Proceedings of the NIPS Workshop on Advances in Structured Learning for Text and Speech Processing*.
- Alexander Clark. 2003. Combining distributional and morphological information for part of speech induction. In *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics - Volume 1*, EACL ’03, pages 59–66, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Jenny R. Finkel and Christopher D. Manning. 2009. Joint parsing and named entity recognition. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, NAACL ’09, pages 326–334, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Jenny Finkel, Shipra Dingare, Huy Nguyen, Malvina Nissim, Christopher Manning, and Gail Sinclair. 2004. Exploiting context for biomedical entity recognition: From syntax to the web. In *In Proceedings of the Joint Workshop on Natural Language Processing in Biomedicine and its applications (JNLPBA-2004)*.
- Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 2005. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, ACL ’05, pages 363–370, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Radu Florian, Abe Ittycheriah, Hongyan Jing, and Tong Zhang. 2003. Named entity recognition through classifier combination. In Walter Daelemans and Miles Osborne, editors, *Proceedings of CoNLL-2003*, pages 168–171. Edmonton, Canada.
- Eduard H. Hovy, Mitchell P. Marcus, Martha Palmer, Lance A. Ramshaw, and Ralph M. Weischedel. 2006. Ontonotes: The 90 In Robert C. Moore, Jeff A. Bilmes, Jennifer Chu-Carroll, and Mark Sanderson, editors, *HLT-NAACL*. The Association for Computational Linguistics.
- Jing Jiang and Chengxiang Zhai. 2006. Exploiting domain structure for named entity recognition. In *In Human Language Technology Conference*, pages 74–81.
- James R. Joel Nothman. 2008. Transforming Wikipedia into Named Entity Training Data. pages 124–132.
- Jun’ichi Kazama and Kentaro Torisawa. 2007. Exploiting Wikipedia as External Knowledge for Named Entity Recognition. In *Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 698–707.
- Jin D. Kim, Tomoko Ohta, Yoshimasa Tsuruoka, Yuka Tateisi, and Nigel Collier. 2004. Introduction to the bio-entity recognition task at JNLPBA. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications*, JNLPBA ’04, pages 70–75, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Dan Klein, Joseph Smarr, Huy Nguyen, and Christopher D. Manning. 2003. Named entity recognition with character-level models. In Walter Daelemans and Miles Osborne, editors, *Proceedings of CoNLL-2003*, pages 180–183. Edmonton, Canada.
- Vijay Krishnan and Christopher D. Manning. 2006. An effective two-stage model for exploiting non-local dependencies in named entity recognition. In *Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*, ACL-44, pages 1121–1128, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Dekang Lin and Xiaoyun Wu. 2009. Phrase Clustering for Discriminative Learning. In *Proceedings of*

- the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1030–1038, Suntec, Singapore, August. Association for Computational Linguistics.
- Dekang Lin, Kenneth Church, Heng Ji, Satoshi Sekine, David Yarowsky, Shane Bergsma, Kailash Patil, Emily Pitler, Rachel Lathbury, Vikram Rao, Kapil Dalwani, and Sushant Narsale. 2010. New Tools for Web-Scale N-grams.
- Andrew McCallum and Wei Li. 2003. Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003 - Volume 4*, CONLL '03, pages 188–191, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Pablo N. Mendes, Max Jakob, Andrés García-Silva, and Christian Bizer. 2011. Dbpedia spotlight: Shedding light on the web of documents. In *Proceedings of the 7th International Conference on Semantic Systems (I-Semantics)*.
- Scott Miller, Jethran Guinness, and Alex Zamanian. 2004. Name tagging with word clusters and discriminative training. In Daniel Marcu Susan Dumais and Salim Roukos, editors, *HLT-NAACL 2004: Main Proceedings*, pages 337–342, Boston, Massachusetts, USA, May 2 - May 7. Association for Computational Linguistics.
- David Nadeau and Satoshi Sekine. 2007. A Survey of Named Entity Recognition and Classification.
- David Nadeau. 2007. *Semi-supervised named entity recognition: learning to recognize 100 entity types with little supervision*. Ph.D. thesis, Ottawa, Ont., Canada, Canada. AAINR49385.
- L. Ratinov and D. Roth. 2009. Design challenges and misconceptions in named entity recognition. In *CoNLL*, 6.
- Sujan Kumar Saha, Sudeshna Sarkar, and Pabitra Mitra. 2009. Feature selection techniques for maximum entropy based biomedical named entity recognition. *J. of Biomedical Informatics*, 42:905–911, October.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: language-independent named entity recognition. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003*, pages 142–147, Morristown, NJ, USA. Association for Computational Linguistics.
- Burr Settles. 2004. Biomedical named entity recognition using conditional random fields and rich feature sets. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications, JNLPBA '04*, pages 104–107, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Valentin I. Spitkovsky, Hiyan Alshawi, Angel X. Chang, and Daniel Jurafsky. 2011. Unsupervised dependency parsing without gold part-of-speech tags. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing (EMNLP 2011)*.
- Charles Sutton and Andrew McCallum. 2006. An Introduction to Conditional Random Fields for Relational Learning. In L. Getoor and B. Taskar, editors, *Introduction to Statistical Relational Learning*. MIT Press.
- Jun Suzuki and Hideki Isozaki. 2008. Semi-Supervised Sequential Labeling and Segmentation Using Giga-Word Scale Unlabeled Data. In *Proceedings of ACL-08: HLT*, pages 665–673, Columbus, Ohio, June. Association for Computational Linguistics.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In Walter Daelemans and Miles Osborne, editors, *Proceedings of CoNLL-2003*, pages 142–147. Edmonton, Canada.
- Maksim Tkatchenko, Alexander Ulanov, and Andrey Simanovsky. 2011. Classifying wikipedia entities into fine-grained classes. In *ICDE Workshops*, pages 212–217.
- Joseph Turian, Lev Ratinov, and Yoshua Bengio. 2010. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, ACL '10*, pages 384–394, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Haochang Wang, Tiejun Zhao, Hongye Tan, and Shu Zhang. 2008. Biomedical named entity recognition based on classifiers ensemble. *IJCSA*, 5(2):1–11.
- Alexander Yeh. 2000. More accurate tests for the statistical significance of result differences.
- GuoDong Zhou and Jian Su. 2004. Exploring deep knowledge resources in biomedical name recognition. In Nigel Collier, Patrick Ruch, and Adeline Nazarenko, editors, *COLING 2004 International Joint workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP) 2004*, pages 99–102, Geneva, Switzerland, August 28th and 29th. COLING.