

# Collective Ontology-based Information Extraction using Probabilistic Graphical Models

Slavko Žitnik and Marko Bajec

University of Ljubljana, Faculty of Computer and Information Science,  
Tržaška cesta 25, SI-1000 Ljubljana  
{slavko.zitnik, marko.bajec}@fri.uni-lj.si

**Abstract.** Traditional information extraction (IE) tasks roughly consist of named-entity recognition, relation extraction and coreference resolution. Much work in this area focuses primarily on separate subtasks where the best performance can be achieved only on specialized domains. In this paper we present a collective IE approach combining all three tasks by employing linear-chain conditional random fields. The usage of probabilistic models enables us to easily communicate between tasks on the fly and error correction during the iterative process execution. We introduce a novel iterative-based IE system architecture with additional semantic and collective feature functions. A new real-world dataset is introduced in the paper against which the proposed system is evaluated. The results are superior over traditional approaches on two tested tasks by error reduction and performance improvements.

**Keywords:** information extraction, named entity recognition, relation extraction, coreference resolution

## 1 Introduction

Machine understanding of textual documents, needed at IE, has been challenging since early computer-era. Russell and Norvig [17] state that IE lies between information retrieval systems, which finds documents related to user's requirements, and text understanding systems that attempt to analyze text and extract their semantic contents. Early IE methods were naive and rule based, then (semi-) automatic approaches of wrapper generation, seed expansion or rule induction were developed and recently machine learning techniques gained popularity. In contrast to standard multi-label and regression classifiers, sequence taggers such as Hidden Markov Models, Maximum Entropy Models and Conditional Random Fields (CRF) have become the most successful. Especially the latter, which supports rich definition of feature functions.

Main IE tasks consist of named entity recognition (e.g. extraction of names, locations, organizations), relation extraction (e.g. identification of connection types among entities) and coreference resolution (e.g. clustering of mentions to an entity). A vast majority of research focuses only on one IE subtask or a pipeline of them not interconnecting them together.

In this paper we propose a collective IE algorithm that iteratively combines all three subtasks. We employ a linear-chain CRF algorithm for every subtask and present additional iterative and semantic features. The use of the same learning techniques enables us to easily use a subtask output labelings as feature inputs for others. In addition to labeling tasks we introduce an entity resolution technique for coreferent mentions matching and merging. Furthermore, during clustering new semantic attributes with coreferent entity values are appended to existing ones, used by CRF's feature functions.

The rest of the paper is organized as follows. Section 2 gives a brief review of related work, focusing mainly on collective IE. Next, a dataset labeling for all three subtasks is presented, followed by CRFs definition, presentation of novel collective algorithm and introduction of new feature functions. Some preliminary results on real-world dataset are discussed in Section 4, and lastly Section 5 concludes the paper and reveals further work.

## 2 Related Work

As mentioned earlier, a majority of research focuses on each subtask separately. The latest research results show the entity extraction problem is quite well solved as best methods achieve 90% F-score on general datasets [14, 3]. On the opposite side, at relation extraction and coreference resolution state-of-the-art methods achieve roughly about 70% F-score [8, 7].

The use of semantic data has also been introduced for IE problems. The general semantic source is Wordnet [10] which contains groups of words with similar meaning - synsets. More exact way of using semantics is an ontology which is defined as explicit specification of conceptualization and provides schema, rules and instances. Recently, ontology-based IE emerged as a subfield of IE [9] because schema, rules and data interoperability can be sufficiently used and merged.

The term collective information extraction was to our knowledge used by Bunescu and Mooney [2] for the first time. They focused only on iterative name entity recognition exploiting mutual influence between possible extractions. Later Nelles and Nazarenko proposed Ontology-based Information extraction [11] that in a cyclic process combines named entity and relation extraction with knowledge integration using an ontology. The proposed system was completely rule based, but it pointed the right direction. The most recent system, Felix [13], was presented by Niu et. al. It is a general IE system based on logical and statistical rules that use Markov Logic Networks. The authors focused on scaling it to large datasets and definition of generally applicable rules. The interesting part is, that their iterative method can combine all three previously mentioned tasks.

Early work in IE was driven by challenges at MUC<sup>1</sup>, CoNLL<sup>2</sup> conferences and ACE<sup>3</sup> program. Along with tasks, the datasets were provided and they are still used for system evaluations.

Conditional Random Fields, a sequence modeling framework, were first presented by Lafferty et. al. [6] and have been since used on various sequence labeling tasks. At proper text labeling and feature induction they were successfully applied to the task of entity recognition [3], relation extraction [8] and coreference resolution [18]. The latter is often accompanied by clustering methods as coreferent mentions need to be merged.

### 3 Collective IE Method

In this section we introduce dataset representation, present used methods and propose a system for collective IE.

#### 3.1 Representation

We treat the tasks of named entity recognition (NER), relation extraction (RE) and coreference resolution (COREF) as sequence labeling tasks.

Let  $\bar{x}^{k_i} = \{x_1^{k_i}, x_2^{k_i}, \dots, x_n^{k_i}\}$  denote the sequence of observable tokens. Index  $k_i$  stands for input words  $w_i$  or additional attributes such as POS tags, phrase boundaries, entity cluster inclusion or pre-calculated values. Each observable sequence is associated with corresponding labeling sequence  $\bar{y}^{l_i}$  where  $l_i \in \{NE, REL, COREF\}$  is defined for named entity, relation and coreference label tags.

We use common IOB notation [16] for all three types of sequence labeling. Tags starting with “B-” denote start of a label type, “I-” the successor of the same type and “O” other types. An example label tag set for person named entities is {B-PER, I-PER, O}. For relations we use labels {B-REL, I-REL, O}. Coreference mentions are labelled using set {B-COREF, I-COREF, O}. Here, tag is labeled as “I-COREF” if and only if it is coreferent with previous (possible distant) tag, labelled as “B-COREF”.

Our problem is now finding the most probable labelings  $\hat{y}^l$  for each of defined subtasks.

#### 3.2 Conditional Random Fields

A Conditional Random Fields (CRFs) [6] are discriminative models and model a single joint distribution  $p(\bar{y}|\bar{x})$  over the predicted sequence  $\bar{y}$  conditioned on  $\bar{x}$ . Observable sequence  $\bar{x}$  typically contains also a number of attributes that can be used when modeling feature functions. Used training labels  $\bar{y}$  relative to position

<sup>1</sup> Message Understanding Conference

<sup>2</sup> Conference on Computational Natural Language Learning

<sup>3</sup> Automatic Content Extraction

$i$  inside feature functions  $f_j$  define the structure of graphical model which can in general be arbitrary.

At CRFs training we are looking for a weight vector  $w$  that assigns best possible labeling  $\hat{y}$  given  $\bar{x}$  for all training examples:

$$\hat{y} = \arg \max_{\bar{y}} p(\bar{y}|\bar{x}; w), \quad (1)$$

using conditional distribution

$$p(\bar{y}|\bar{x}; w) = \frac{\exp(\sum_{j=1}^J w_j \sum_{i=1}^n f_j(\bar{y}, \bar{x}, i))}{Z(\bar{x}, w)} \quad (2)$$

Vector  $w$  contains a real number for every possible input ( $J$  inputs) to  $f_j$ . ( $Z(\bar{x}, w)$  is a normalization constant over all possible labelings of  $\bar{y}$ ). When distance between two addressing labels inside feature functions is long, exact inference is intractable due to exponential number of partial sequences and thus approximate algorithms must be used. We therefore use feature functions that depend only on single label ( $y_i$ ) and two consecutive labels ( $y_{i-1}, y_i$ ). This type of CRF is also known as linear chain CRFs (LCCRFs) which underlying graphical structure forms a chain and have been rather successful in IE tasks. Using LCCRFs, training and inference can be easily solved using forward-backward method and Viterbi algorithm [15].

### 3.3 Collective approach

We propose a collective IE algorithm combining tasks of entity recognition, relation extraction and coreference resolution. A high level implementation of iterative training and labeling algorithm is shown as Algorithm (1) and Algorithm (2).

The input parameters for training Algorithm (1) are sentences, tokenized by words with additional attributes  $\bar{x}^k$ , true named entities, relations and coreferences labelings denoted as  $\bar{y}^l$  and number of maximum possible iterations. The final training result is a 3-tuple of trained classifiers for each task. These classifiers can be used independently, but are trained to be used by Algorithm (2) to get best results.

During each learning iteration in Algorithm (1), feature function vectors are initialized and then classifiers are independently learned. After that we update/create additional attributes for next iteration of training. At that step we perform collective entity resolution using attribute, relational and semantic similarity measures as proposed and evaluated in [20, 21]. As a result we get clusters of coreferent entities which values are used at initialization of feature vectors in next iteration. Similarly we provide additional attributes by tagging the input sequence using latest classifiers. For example, let have an input sequence: “*John has left ACME. ... When he worked at ACME as a student.*” In iteration  $i$  ACME was not recognized as a company in first sentence, but all others were correctly labeled and pronoun *he* was merged with *John*. In the next iteration feature functions can use distant relation *workedAt* **Company** which results in correct first

---

**Algorithm 1** Collective IE Training
 

---

**Input:**  $\bar{x}^k, \bar{y}^l, \text{maxIter}$   
**Output:** classifiers (cNE, cREL, cCOREF)  
 1: Initialize coref. clusters as  $C = \emptyset$   
 2:  $i \leftarrow 0$   
 3: **while**  $i < \text{maxIter}$  **and**  $\text{prevScoreDiff}() < \varepsilon$  **do**  
   4: Initialize feature functions  
   5:  $\text{cNE} \leftarrow \text{LCCRF}(\bar{x}^k, \bar{y}^{NE})$   
   6:  $\text{cREL} \leftarrow \text{LCCRF}(\bar{x}^k, \bar{y}^{REL})$   
   7:  $\text{cCOREF} \leftarrow \text{LCCRF}(\bar{x}^k, \bar{y}^{COREF})$   
   8:  $C \leftarrow \text{entityResolution}(\bar{x}^k, \bar{y}^l)$   
   9:  $\bar{x}^{I-NE} \leftarrow \text{cNE.tag}(\bar{x}^k)$   
   10:  $\bar{x}^{I-REL} \leftarrow \text{cREL.tag}(\bar{x}^k)$   
   11:  $\bar{x}^{I-COREF} \leftarrow \text{cCOREF.tag}(\bar{x}^k)$   
   12:  $i \leftarrow i + 1$   
 13: **end while**  
 14: **return** (cNE, cREL, cCOREF)

---

sentence labeling. Iterating ends when classifier’s labelings over iterations converge or maximum number of iterations is achieved. We will empirically define  $\text{maxIter}$  and  $\text{prevScoreDiff}$  in further work with all three subtasks.

Algorithm (2) introduces iterative labeling and is very similar to training algorithm. Only feature initialization, tagging and coreference clustering is used until there are no labeling differences over two sequential iterations or maximum number of them is reached.

---

**Algorithm 2** Collective IE Labeling
 

---

**Input:**  $\bar{x}^k, (\text{cNE}, \text{cREL}, \text{cCOREF}), \text{maxIter}$   
**Output:** labelings and coreference clusters  
 1: Initialize coref. clusters as  $C = \emptyset$   
 2:  $i \leftarrow 0$   
 3: **while**  $\text{labelingsChanged}()$  **and**  $i < \text{maxIter}$  **do**  
   4: Initialize feature functions  
   5:  $\bar{x}^{I-NE} \leftarrow \text{cNE.tag}(\bar{x}^k)$   
   6:  $\bar{x}^{I-REL} \leftarrow \text{cREL.tag}(\bar{x}^k)$   
   7:  $\bar{x}^{I-COREF} \leftarrow \text{cCOREF.tag}(\bar{x}^k)$   
   8:  $C \leftarrow \text{entityResolution}(\bar{x}^k, \bar{y}^l)$   
   9:  $i \leftarrow i + 1$   
 10: **end while**  
 11: **return** ( $\bar{x}^{NE}, \bar{x}^{REL}, C$ )

---

### 3.4 Features

The selection of feature functions is an essential step for successfully training CRF classifiers.

We use proposed named entity feature functions by Manning et. al. [3], relation-specific features proposed by Li et. al. [8] and coreference-specific features by McCallum et. al. [18] and Ng. and Cardie [12]. The union of all features across tasks represents word, text preprocessing (i.e. POS tags, lemmas, Parse trees) and word shape features.

In Table 1 we introduce additional iterative and semantic feature functions. It is worth mentioning that some local attributes are equivalent to long-distance that can be modeled as arbitrary structured CRF and are here a result of entity resolution.

**Table 1.** Linear-chain feature function templates.  $i$  indicates current position and  $j$  offset relative to  $i$ . Functions depending only at one label generate  $\#labels$  features and  $\#labels^2$  features for depending on two consecutive labels.

Feature Function Description	Feature Template	Example
single cluster relation	$c_i(-1 \leq i \leq 1)$	$c_i$ is <i>works at</i>
single cluster entity tag	$c_i(-1 \leq i \leq 1)$	$c_i$ is <i>I-PER</i>
single cluster entity name	$c_i(-1 \leq i \leq 1)$	$c_i$ is <i>Chuck Norris</i>
single previous iter NE tag	$n_i(-2 \leq i \leq 2)$	$n_i$ is <i>B-PER</i>
single previous iter REL tag	$r_i(-2 \leq i \leq 2)$	$r_i$ is <i>I-REL</i>
single previous iter COREF tag	$co_i(-2 \leq i \leq 2)$	$co_i$ is <i>B-COREF</i>
two entity tags at coreferences	$c_{i+j-1} \& c_{i+j}(-1 \leq j \leq 2)$	$c_{i+j-1}$ is <i>I-ORG</i> and $c_{i+j}$ is <i>B-PER</i>
two iter NE tags	$n_{i+j-1} \& n_{i+j}(-1 \leq j \leq 2)$	$n_{i+j-1}$ is <i>O</i> and $n_{i+j}$ is <i>B-ORG</i>
two iter REL tags	$r_{i+j-1} \& r_{i+j}(-1 \leq j \leq 2)$	$r_{i+j-1}$ is <i>B-REL</i> and $r_{i+j}$ is <i>I-REL</i>
two iter COREF tags	$co_{i+j-1} \& co_{i+j}(-1 \leq j \leq 2)$	$co_{i+j-1}$ is <i>O</i> and $co_{i+j}$ is <i>B-COREF</i>

## 4 Evaluation and Results

We conducted some preliminary analysis of part of proposed method. We employed only tasks of named entity recognition and relation extraction to show the results are promising and it is worth to build the whole system.

We tested methods on real-world news dataset in Slovene language which is publicly available <sup>4</sup>. The topic of the articles is mainly political. It contains 6034 word tokens within 285 sentences. The tokens are annotated according to

<sup>4</sup> [http://zitnik.si/mediawiki/index.php?title=File:Rtvslo\\_dec2011.tsv](http://zitnik.si/mediawiki/index.php?title=File:Rtvslo_dec2011.tsv)

proposed representation in Section 3.1 (their distribution is shown in Table 2). We additionally lemmatized and part-of-speech (POS) tagged the whole corpus using slovene POS tagger [4]. Slovene is morphologically complex language and therefore around thousand different POS tags exist [5]. Within the dataset there are 315 distinct POS tags labeled.

**Table 2.** A distribution of BIO tags following proposed representation in Section 3.1

Type	B-	I-	O
Named entities	293	233	5508
Relations	32	24	5978
Coreferences	274	249	5511

For evaluation purposes we implemented linear-chain CRFs using stochastic gradient ascent learning algorithm and Viterbi for decoding [1] with rich feature function API that supports arbitrary implementation. Whole implementation with additional broader framework features is available online <sup>5</sup>.

In Table 3 we show achieved MAF (macro averaged F-score) measure [19] when training independently, using traditional “pipeline” approach and by employing part of proposed collective algorithm. In macro-averaging, F-measure is computed locally over each category first and then the average over all categories is taken. All approaches use feature functions leveraging learning label, upper case of one and two consecutive words, prefixes and suffixes of length two and three and POS, words and lemma features. At collective approach, additional iterative feature functions for named entity recognition and relation extraction are used. Results show the collective approach outperforms other two in both labeling tasks.

**Table 3.** Comparison of macro averaged F-score on real world dataset by independent, pipeline and collective learning

	Entity Recognition	Relation Extraction
Independent	.57	.71
Pipeline	.57	.72
Collective	<b>.59</b>	<b>.74</b>

## 5 Conclusions and Future Work

Paper proposes a collective information extraction algorithm, which combines tasks of named entity recognition, relation extraction and coreference resolu-

<sup>5</sup> <https://bitbucket.org/szitnik/iobie/>

tion. We introduce iterative training and labeling algorithm, present new iterative feature functions and show preliminary experimental results which show improvements over traditional approaches.

Future work will include implementation of the whole proposed iterative system, modeling of new features and introduction of parallelization algorithms. After that we will be able to compare results to others on larger datasets. Ontologies will be incorporated into the algorithm to use additional manual patterns, constraints and to better connect with entity resolution module.

## 6 Acknowledgments

The work has been supported by the Slovene Research Agency ARRS within the research program P2-0359 and part financed by the European Union, European Social Fund.

## References

1. Léon Bottou. Large-scale machine learning with stochastic gradient descent. In Yves Lechevallier and Gilbert Saporta, editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010)*, pages 177–187, Paris, France, August 2010. Springer.
2. Razvan Bunescu and Raymond J. Mooney. Collective information extraction with relational markov networks. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, ACL '04, Stroudsburg, PA, USA, 2004. Association for Computational Linguistics.
3. Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 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, 2005. Association for Computational Linguistics.
4. Miha Grčar, Jan Rupnik, Matjaž Juršič, and Simon Krek. Slovene pos tagger. <http://označevalnik.slovensščina.eu/Vsebina/Sl/ProgramskaOprema/Meta.aspx>.
5. Primož Jakopin and Aleksandra Bizjak. Part-of-speech tagging of slovenian text. *Slavistična revija*, 45(3-4):513–532, 1997.
6. John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, ICML '01, pages 282–289, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
7. Heeyoung Lee, Yves Peirsman, Angel Chang, Nathanael Chambers, Mihai Surdeanu, and Dan Jurafsky. Stanford's multi-pass sieve coreference resolution system at the conll-2011 shared task. In *Proceedings of the CoNLL-2011 Shared Task*, 2011.
8. Y. Li, J. Jiang, H.L. Chieu, and K.M.A. Chai. Extracting relation descriptors with conditional random fields. pages 392–400, Thailand, 2011. Asian Federation of Natural Language Processing.
9. L. McDowell and M. Cafarella. Ontology-driven information extraction with ontosyphon. *The Semantic Web-ISWC 2006*, page 428–444, 2006.

10. George A. Miller. Wordnet: A lexical database for english. *Communications of the ACM*, 38:39–41, 1995.
11. Claire Nedellec and Adeline Nazarenko. Ontologies and information extraction. *CoRR*, abs/cs/0609137, 2006.
12. Vincent Ng and Claire Gardent. Improving machine learning approaches to coreference resolution. In *ACL*, pages 104–111, 2002.
13. Feng Niu, Ce Zhang, Christopher Ré, and Jude W. Shavlik. Felix: Scaling inference for markov logic with an operator-based approach. *CoRR*, abs/1108.0294, 2011.
14. Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R. Curran. Learning multilingual named entity recognition from wikipedia. *Artificial Intelligence*, (0):–, 2012.
15. L. R. Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
16. L.A. Ramshaw and M.P. Marcus. Text chunking using transformation-based learning. In *Proceedings of the Third ACL Workshop on Very Large Corpora*, page 82–94, 1995.
17. Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 2003.
18. Michael L. Wick, Aron Culotta, Khashayar Rohanimanesh, and Andrew McCallum. An entity based model for coreference resolution. In *SDM*, pages 365–376, 2009.
19. A. Özgür, L. Özgür, and T. Güngör. Text categorization with class-based and corpus-based keyword selection. *Computer and Information Sciences-ISCIS 2005*, page 606–615, 2005.
20. Lovro Šubelj, David Jelenc, Eva Zupančič, Dejan Lavbič, Denis Trček, Marjan Krisper, and Marko Bajec. Merging data sources based on semantics, contexts and trust. In *The IPSI BgD Transactions on Internet Research*, volume 7, page 18, 2011.
21. Slavko Žitnik, Lovro Šubelj, Dejan Lavbič, Olegas Vasilecas, and Marko Bajec. Contextual data matching and merging using semantics, trust and ontologies. *Informatica - in review*, 2012.