

Definition Extraction with Balanced Random Forests

No Author Given

No Institute Given

Abstract. In this paper we propose a machine learning approach to the task of identifying definitions in Polish documents. Specifics of the problem domain and characteristics of the available dataset have been taken into consideration, by carefully choosing and adapting a classification method to highly imbalanced and noisy data. We evaluate the performance of a Random Forest-based classifier in extracting definitional sentences from natural language text and give a comparison with previous work.

1 Introduction

Natural Language Processing (NLP) tasks often involve heavily imbalanced data, with a dominating “uninteresting” class and a minority “interesting” class. One such task is that of definition extraction, where a set of sentences is to be classified into definitional and non-definitional sentences. There may be as many as 20 non-definition sentences for any single definition sentence in an instructive text, but it is the latter class that a definition extraction system is interested in.

The usual Machine Learning (ML) classifiers, ranging from naïve bayesian methods, through decision trees, perceptrons and various lazy learners, to the currently very popular classifiers based on Support Vector Machines (SVMs) and on Adaboost, do not work well in such cases, even when trained with subsampling (of uninteresting examples) or oversampling (of the interesting examples). The problem is that such classifiers attempt to minimise the overall error rate, rather than concentrating on the interesting class. In case of a dataset with a 1:20 ratio of interesting to uninteresting cases, it is difficult to beat a classifier uniformly assigning each new item to the uninteresting class: such a classifier reaches the overall accuracy higher than 95%, but at the cost of misclassifying all the interesting cases!

The problem of heavily imbalanced data has already been addressed in the ML community, where most solutions consist either in the assignment of a high cost to the misclassification of the minority class or in subsampling and/or oversampling. A novel approach to the problem has been proposed in Chen *et al.* 2004 and it consists in a modification of the Random Forest (Breiman, 2001) classifier.

Random Forest (RF) is a homogeneous ensemble of regression (e.g., CART) or decision (e.g., C4.5) unpruned trees, where — at each node of the tree —

a subset of all attributes is randomly selected and the best attribute on which to further grow the tree is taken from that random set. Additionally, Random Forest is an example of the bagging (bootstrap aggregating) method, i.e., each tree is trained on a set bootstrapped¹ from the original training set. Decisions are reached by simple voting.

Balanced Random Forest (BRF; Chen *et al.* 2004) is a modification of RF, where for each tree two bootstrapped sets of the same size, equal to the size of the minority class, are constructed: one for the minority class, the other for the majority class. Jointly, these two sets constitute the training set.

The aim of this paper is to demonstrate that BRF is a technique well-suited to the difficult problem of definition extraction and, by extension, other NLP tasks. When trained on the dataset of Polish instructive texts introduced in Przepiórkowski *et al.* 2007b,a, BRF-based classifiers give better results than manual definition extraction grammars (Przepiórkowski *et al.*, 2007a) or the usual ML classifiers, even when combined with some *a priori* linguistic knowledge (Degórski *et al.*, 2008).

In what follows we first introduce the attribute space assumed here (§2), then describe the used classification approach (§3) and present the results of our experiments (§4). Finally, we outline work conducted previously in the field (§5) and conclude with possibilities of further research (§6).

2 Feature Selection

Employing any machine learning approach to unstructured data requires that data is represented in the form of feature values, either binary, numeric or nominal. We use a relatively straightforward approach of n -gram representation of the sentences in the available document set. Each sentence is represented by a vector of binary values, where each value indicates whether a particular n -gram is present in the corresponding sentence. The n -grams consist of base forms of words, their parts of speech and grammatical cases that appear in the greatest number of sentences in all documents. We individually count the occurrences of each of the n -grams in sentences marked as definitions and non-definitions. Both lists are then combined and a number of most common entries is selected to form a dictionary of features used for sentence description.

Even by limiting the length of generated n -grams to $n \leq 3$ and having a choice of three distinct n -gram types: base word form (further denoted as *base*), part of speech of the word (*ctag*) and its grammatical case (*case*), we face a problem of many possible dictionary configurations, selecting from the set of $3^1 + 3^2 + 3^3 = 39$ possibilities. Including too many n -gram types would result in an extremely large attribute space, while including too few in reducing the potential classification accuracy. We approached the problem by measuring the average value of the χ^2 statistic of each of the possible n -gram types with respect

¹ That is, examples in such a bootstrapped training set are uniformly and randomly drawn *with replacement* from the original training set. As a result, some examples will be repeated while other will not make it to the bootstrapped set.

to the class attribute. This was performed on a training set consisting of all the available documents, on the basis of 100 n -grams for each of the 39 types. Table 1 presents a list of the 20 n -gram types with the highest average χ^2 value.

Table 1: Top 20 values of the χ^2 statistic of possible n -gram permutations.

rank	n -gram	average χ^2	rank	n -gram	average χ^2
1	<i>base</i>	21.04	11	<i>base base ctag</i>	16.70
2	<i>ctag ctag case</i>	18.91	12	<i>ctag base ctag</i>	16.29
3	<i>ctag base</i>	18.53	13	<i>ctag ctag base</i>	14.77
4	<i>base case</i>	18.45	14	<i>ctag case</i>	14.69
5	<i>base ctag</i>	17.92	15	<i>ctag ctag ctag</i>	14.63
6	<i>base base</i>	17.81	16	<i>base ctag case</i>	14.52
7	<i>base base case</i>	17.73	17	<i>base base base</i>	14.33
8	<i>ctag base case</i>	17.43	18	<i>ctag</i>	13.88
9	<i>ctag ctag</i>	17.11	19	<i>ctag case ctag</i>	13.65
10	<i>ctag base base</i>	16.73	20	<i>base ctag ctag</i>	13.59

Unfortunately, just taking a number of attributes from the top of this list does not guarantee the best possible selection of n -gram types. This is because certain attribute pairs may be statistically dependent and introducing both of them into the dictionary would result in noise, instead of meaningful data for the classifier. Having experimented with different attribute configurations, we have chosen the following heuristic procedure of attribute selection: we take one attribute at a time from the sorted list, starting from the top, and reject these n -grams of length $n = 3$, for which another trigram with one of the same feature types has already been selected. The resulting set of 10 selected n -gram types is presented in Table 2.

Table 2: The selected set of n -gram types.

no.	n -gram	no.	n -gram
1	<i>base</i>	6	<i>base base</i>
2	<i>ctag ctag case</i>	7	<i>ctag ctag</i>
3	<i>ctag base</i>	8	<i>ctag case</i>
4	<i>base case</i>	9	<i>base base base</i>
5	<i>base ctag</i>	10	<i>ctag</i>

For comparison purposes, we also present here the results of experiments on a dataset from the work of Degórski *et al.* 2008, where the set of n -grams has been selected in a different manner and using a different set of attributes. Specifically, this dataset, referenced later as the “baseline dataset”, has been created by using

the 100 most common uniform unigrams, bigrams and trigrams of base forms, parts of speech and cases (i.e., *base*, *base-base*, . . . , *ctag-ctag-ctag*).

3 Classifying Imbalanced Data

As noted earlier, the available dataset of definitional and non-definitional sentences is highly imbalanced and consists of 10830 sentences, 546 of which contain — or are a part of — definitions. Consequently, any successful classification-based approach to extraction of definitions from this data must take into consideration — either explicitly or implicitly — the difference in training samples from both categories.

The most common way of dealing with imbalanced data is introducing appropriately weighted costs for specific classes or sampling the available training set. Balanced Random Forest is an approach where equalizing the influences of classes is not performed externally to classification algorithm by evaluating weights, but is integrated in the very process. Here, for the task of extracting definitions from a set of documents by sentence classification, we use the following algorithm, based on Chen *et al.* 2004:

- split the training corpus into definitions and non-definitions; let us assume that there are n_d definitions and n_{nd} non-definitions;
- construct k trees, each in the following way:
 - draw a bootstrap sample of size n_d of definitions, and a bootstrap sample of the same size n_d of non-definitions;
 - learn the tree (without pruning) using the CART algorithm, on the basis of the sum of the two bootstrap samples as the training corpus, but:
 - at each node, first select at random m features (variables) from the set of all M features ($m < M$; selection without replacement), and only then select the best feature (out of these m features) for this node; this random selection of m features is repeated for each node;
- the final classifier is the ensemble of the k trees and decisions are reached by simple voting.

We have chosen the value of m to be equal to \sqrt{M} in all the experiments.

As Random Forest is a well known classifier and widely covered in the literature, it also allows having a greater insight into the results produced by the BRF approach. RFs have been verified to be suitable both for large and highly dimensional data, as is the case in natural language processing. They also provide means of estimating the classification error rate without performing a full cross-validation procedure and for estimating variable importance and variable interactions. In our current experiments we have not performed such estimations, as we are more interested in selecting the optimal set of n -gram types, than comparing the importance of particular features.

4 Experimental Results

We use several statistical parameters to describe and compare the results of the proposed classification approach: recall and precision are the most commonly calculated information retrieval performance measures. We assume the sentences marked as definitions to be the set of relevant documents in the retrieval task:

$$precision = \frac{|\{\text{definitions}\} \cap \{\text{retrieved sentences}\}|}{|\{\text{retrieved sentences}\}|} \quad (1)$$

$$recall = \frac{|\{\text{definitions}\} \cap \{\text{retrieved sentences}\}|}{|\{\text{definitions}\}|} \quad (2)$$

For a single-valued performance indicator, we use the F-measure, both in the form used in the previous papers on Polish definition extraction (marked as F_α) and in the more common sense (marked as F_β). For F_1 we just use F_1 (as $F_{\alpha=1} = F_{\beta=1}$):

$$F_\alpha = \frac{(1 + \alpha) \cdot precision \cdot recall}{\alpha \cdot precision + recall} \quad (3)$$

$$F_\beta = \frac{(1 + \beta^2) \cdot precision \cdot recall}{\beta^2 \cdot precision + recall} \quad (4)$$

Finally, we also calculate the area under the ROC curve (AUC), which is another single-valued measure of retrieval accuracy, but not tied to a single probability threshold value, like the F-measure. Still, because in the task of definition extraction we are more interested in maximizing the recall value (in other words: minimizing the false negative rate), we compare all further experiment results on the basis of $F_{\alpha=2}$ and $F_{\beta=2}$ values.

Our initial experiments aimed at verifying whether any additional preprocessing of the available data, commonly applied to text classification problems, would result in improving the accuracy of definition extraction. Firstly, we have included the information about the relative position of an n -gram in a sentence into the feature vector. By dividing the sentences into three equal parts and counting the n -gram occurrences in each of the parts separately, we have increased the attribute space three times, but achieved no increase in performance (Table 3). We may speculate that the positional information introduced too much noise, as the available dataset was too small to benefit from the significantly larger feature space.

Similarly, there was no gain in definition extraction accuracy after including the information about the actual number of occurrence counts of particular n -grams in the analyzed sentences. This may also be explained by a relatively small size of the available dataset and sparseness of the feature vector. The calculated numbers of occurrences were negligibly small and provided no additional information to the classifier.

Finally, applying a stop-list of most common words and filtering non-alphanumeric characters from the documents also proved to reduce both the value of $F_{\alpha=2}$ and $F_{\beta=2}$ measures. Thus, neither of the attribute modifications and data

Table 3: The influence of additional preprocessing steps on classification accuracy. Ten-fold cross-validation results, with 100 iterations of random trees generation.

dataset	precision	recall	F_1	$F_{\alpha=2}$	$F_{\alpha=5}$	$F_{\beta=2}$	$F_{\beta=5}$	AUC
base	18.11%	66.10%	28.43%	35.10%	45.85%	43.20%	59.99%	82.36%
n -gram position	16.20%	63.90%	25.85%	32.25%	42.86%	40.22%	57.40%	81.20%
n -gram occurrence	17.20%	65.00%	27.20%	33.74%	44.42%	41.78%	58.72%	81.40%
base form stoplist	17.30%	63.00%	27.15%	33.50%	43.74%	41.22%	57.19%	81.40%

preprocessing steps mentioned above have been used in further experiments. A detailed comparison of each of the approaches has been presented in Figure 2a.

In an effort to determine the optimal size of feature space for classification, we have conducted a series of experiments with an increasing number of n -grams used for sentence representation (Table 4 and Figure 1a). On the basis of the results, we have decided to use 100 n -grams of each type in further experiments, as increasing their number above that threshold does not seem to have any positive influence on the classification accuracy. By choosing that number, we obtained a training set consisting of 10830 instances and 929 attributes (as there are less than 100 different n -grams of the type *ctag*).

Table 4: The influence of the number of used n -grams of each type on classification accuracy. Ten-fold cross-validation results, with 100 iterations of random trees generation.

n -grams	precision	recall	F_1	$F_{\alpha=2}$	$F_{\alpha=5}$	$F_{\beta=2}$	$F_{\beta=5}$	AUC
10	14.41%	57.69%	23.06%	28.83%	38.45%	36.04%	51.72%	76.64%
20	17.20%	63.71%	27.09%	33.51%	43.92%	41.35%	57.71%	81.65%
30	18.66%	65.37%	29.03%	35.64%	46.13%	43.56%	59.63%	82.74%
40	19.05%	66.84%	29.65%	36.40%	47.13%	44.51%	60.96%	82.58%
50	19.33%	67.94%	30.10%	36.96%	47.87%	45.20%	61.95%	83.12%
60	19.25%	67.22%	29.93%	36.72%	47.49%	44.86%	61.34%	82.89%
70	19.14%	66.84%	29.76%	36.51%	47.22%	44.61%	60.99%	83.23%
80	19.66%	67.20%	30.42%	37.21%	47.90%	45.29%	61.48%	83.72%
90	19.78%	69.42%	30.79%	37.80%	48.95%	46.22%	63.31%	84.48%
100	20.10%	70.10%	31.24%	38.32%	49.55%	46.81%	63.98%	83.80%
110	19.60%	68.10%	30.44%	37.32%	48.22%	45.55%	62.18%	84.10%
120	19.60%	67.80%	30.41%	37.26%	48.09%	45.45%	61.94%	84.10%

As the accuracy of Random Forest classification depends heavily on the number of generated random trees used in voting, we have conducted the experiments both on the current dataset and on the baseline dataset provided by Degórski *et al.* 2008 for several different numbers of iterations (Tables 5 and 6, Figure 1b). We have performed ten-fold cross-validation experiments instead of counting the

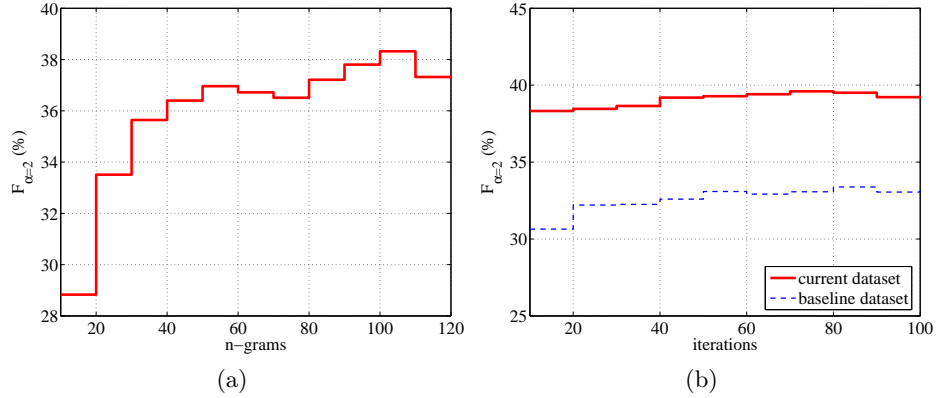


Fig. 1: (a) Performance of classification with respect to the number of used n -grams, (b) a comparison between classification performance using the baseline dataset and the current dataset for different number of iterations.

out-of-bag error of the bagging classifier, so as to make the results as closely comparable with those of Degórski *et al.* 2008 as possible. The detailed comparison of both sets, with respect to BRF classification accuracy for the number of iterations which proved to give the best results for each of the sets, is presented in Figure 2b.

Table 5: Ten-fold cross-validation results of the baseline dataset classification.

iterations	precision	recall	F_1	$F_{\alpha=2}$	$F_{\alpha=5}$	$F_{\beta=2}$	$F_{\beta=5}$	AUC
100	15.11%	63.03%	24.38%	30.64%	41.23%	38.57%	56.18%	80.55%
200	16.12%	64.32%	25.78%	32.21%	42.93%	40.25%	57.69%	81.26%
300	16.22%	63.77%	25.86%	32.25%	42.84%	40.20%	57.31%	81.57%
400	16.50%	63.58%	26.20%	32.59%	43.09%	40.48%	57.29%	81.62%
500	16.81%	64.13%	26.64%	33.09%	43.65%	41.03%	57.87%	81.70%
600	16.71%	63.94%	26.50%	32.92%	43.46%	40.85%	57.67%	81.70%
700	16.87%	63.59%	26.67%	33.07%	43.51%	40.92%	57.47%	81.83%
800	17.04%	64.13%	26.93%	33.38%	43.91%	41.30%	57.97%	81.87%
900	16.86%	63.59%	26.65%	33.05%	43.50%	40.91%	57.46%	81.87%
1000	16.91%	63.96%	26.75%	33.18%	43.70%	41.09%	57.78%	81.89%

As may be seen from the results of the consecutive experiments, increasing the number of generated random trees improves the accuracy of definitional sentences classification only up to a certain point. Above that threshold the performance reaches a plateau and no further iterations are necessary.

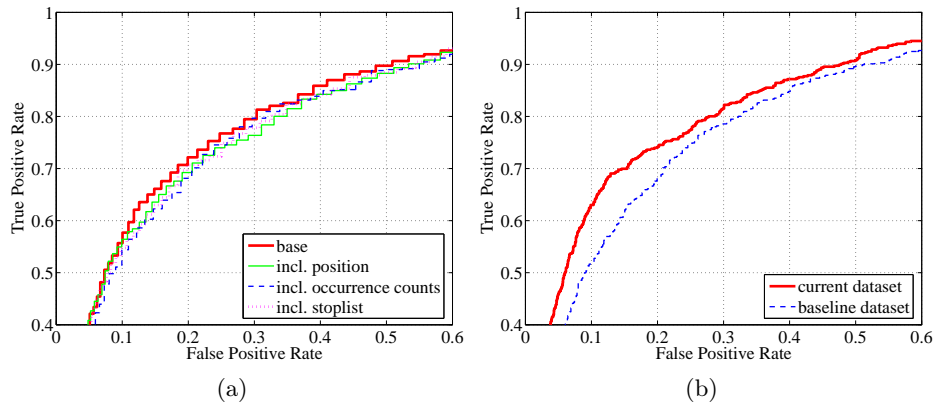


Fig. 2: ROC curve of classification: (a) using additional data preprocessing steps, (b) using the baseline dataset and the current dataset.

Table 6: Ten-fold cross-validation results of the current dataset classification.

iterations	precision	recall	F_1	$F_{\alpha=2}$	$F_{\alpha=5}$	$F_{\beta=2}$	$F_{\beta=5}$	AUC
100	20.10%	70.10%	31.24%	38.32%	49.55%	46.81%	63.98%	83.80%
200	20.46%	68.67%	31.53%	38.46%	49.31%	46.67%	62.96%	84.35%
300	20.62%	68.68%	31.72%	38.65%	49.46%	46.84%	63.03%	84.49%
400	20.98%	69.22%	32.20%	39.19%	50.04%	47.42%	63.60%	84.59%
500	21.13%	68.86%	32.34%	39.28%	50.03%	47.43%	63.36%	84.71%
600	21.24%	68.86%	32.47%	39.41%	50.13%	47.54%	63.39%	84.73%
700	21.37%	69.04%	32.64%	39.60%	50.33%	47.74%	63.58%	84.72%
800	21.29%	69.04%	32.54%	39.51%	50.25%	47.66%	63.56%	84.78%
900	21.11%	68.68%	32.29%	39.22%	49.93%	47.34%	63.20%	84.78%
1000	21.20%	68.68%	32.40%	39.32%	50.01%	47.43%	63.23%	84.79%

While the use of Balanced Random Forest classification method alone significantly improves the definition extraction performance over other pure machine learning based approaches (e.g., as reported by Degórski *et al.* 2008), it is worth pointing out that a careful feature selection is an equally important step. We achieve an over 18% increase in accuracy, as indicated by the $F_{\alpha=2}$ measure, by describing the sentences with a more representative set of attribute types.

5 Previous Work

To the best of our (and Google’s) knowledge, there is no previous NLP work taking advantage of the Balanced variety of RFs. Apparently, the first NLP applications of the plain Random Forests are those reported in Nielsen and Pradhan 2004, for PropBank-style (Kingsbury and Palmer, 2002) role classification, and

in Xu and Jelinek 2004 (followed by a series of papers by the same authors, culminating in Xu and Jelinek 2007), where they are used in the classical language modelling task (predicting a sequence of words) for speech recognition and give better results than the usual n -gram based approaches.

On the other hand, there is some substantial previous work on definition extraction, as this is a subtask of many applications, including terminology extraction (Pearson, 1996), the automatic creation of glossaries (Klavans and Muresan, 2000, 2001), question answering (Miliaraki and Androutsopoulos, 2004; Fahmi and Bouma, 2006), learning lexical semantic relations (Malaisé *et al.*, 2004; Storrer and Wellinghoff, 2006) and the automatic construction of ontologies (Walter and Pinkal, 2006). Despite the current dominance of the ML paradigm in NLP, tools for definition extraction are invariably language-specific and involve shallow or deep processing, with most work done for English (Pearson, 1996; Klavans and Muresan, 2000, 2001) and other Germanic languages (Fahmi and Bouma, 2006; Storrer and Wellinghoff, 2006; Walter and Pinkal, 2006), as well as French (Malaisé *et al.*, 2004).

When ML methods are used, it is in combination with linguistic processing. For example, Fahmi and Bouma 2006 applied a robust wide-coverage parser of Dutch to select candidate definition sentences, which were then subject to a ML classifier. They experimented with three classifiers (Naïve Bayes, SVM and Maximum Entropy) and a number of possible feature configurations and obtained the best results for the Maximum Entropy classifier and feature configurations, which included some syntactic features.

For Polish, first attempts at constructing definition extraction systems are described — in the context of other Slavic languages — in Przepiórkowski *et al.* 2007b, and improved results are presented in Przepiórkowski *et al.* 2007a. In that work definitions were identified on the basis of a manually constructed partial grammar (a cascade of regular grammars over morphosyntactically annotated XML-encoded texts), with the best grammar giving the precision of 18.69 and recall of 59.34, which amounts to $F_{\alpha=2} = 34.39$. Przepiórkowski *et al.* 2007a note that these relatively low results are at least partially due to the inherent difficulty of the task: the inter-annotator agreement measured as Cohen’s κ is only 0.31 (the value of 1 would indicate perfect agreement, the value of 0 — complete randomness). The same dataset was used in the experiments reported here.

An approach more directly comparable to ours is presented in Degórski *et al.* 2008. The general idea is analogous to that of Fahmi and Bouma 2006: first candidate definition sentences are selected via linguistic methods and then they are classified using ML methods. What is novel in Degórski *et al.* 2008 is the very basic character of the linguistic knowledge (a small low-precision collection of n -grams typical for definitions, including the copula, sequences corresponding to *that is* and *i.e.*, etc.), and the use of ensembles of classifiers in the second stage. The best results reported there, the precision of 19.94, recall of 69.23, and $F_{\alpha=2} = 37.95$, are significantly better than those of Przepiórkowski *et al.* 2007a,

but still, despite some use of *a priori* language-specific knowledge, worse than the pure ML results reported here.

6 Conclusions and Future Work

Our currently reported results seem to restore hope in the machine learning approach to the vaguely specified task of definition extraction from a small set of text documents. It is usually the case that the smaller, less structured and more noisy the available training data, the lesser is the advantage of such methods over hand-crafted rules and grammars, utilizing linguistic knowledge. Thus, achieving better results in such circumstances by a pure machine learning approach seems to justify the necessary work on feature and classification method selection.

It would still be interesting to combine the current classification method with manually constructed grammars, similarly as in Degórski *et al.* 2008, to see if such a sequential processing scheme would further improve the definition extraction performance. On the basis of the experiments described there, we might expect a considerable increase in retrieval precision, at the cost of a slight decrease in recall.

References

- Breiman, L. (2001). Random forests. *Machine Learning*, **45**, 5–32.
- Chen, C., Liaw, A., and Breiman, L. (2004). Using random forest to learn imbalanced data. Technical Report 666, University of California, Berkeley. <http://www.stat.berkeley.edu/tech-reports/666.pdf>.
- Degórski, Ł., Marcińczuk, M., and Przepiórkowski, A. (2008). Definition extraction using a sequential combination of baseline grammars and machine learning classifiers. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation, LREC 2008*, Marrakech. ELRA. Forthcoming.
- Fahmi, I. and Bouma, G. (2006). Learning to identify definitions using syntactic features. In *Proceedings of the EACL 2006 workshop on Learning Structured Information in Natural Language Applications*.
- Kingsbury, P. and Palmer, M. (2002). From TreeBank to PropBank. In *Proceedings of the Third International Conference on Language Resources and Evaluation, LREC 2002*, pages 1989–1993, Las Palmas. ELRA.
- Klavans, J. L. and Muresan, S. (2000). DEFINDER: Rule-based methods for the extraction of medical terminology and their associated definitions from on-line text. In *Proceedings of the Annual Fall Symposium of the American Medical Informatics Association*.
- Klavans, J. L. and Muresan, S. (2001). Evaluation of the DEFINDER system for fully automatic glossary construction. In *Proceedings of AMIA Symposium 2001*.
- Lin, D. and Wu, D., editors (2004). *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP 2004)*, Barcelona. ACL.

- Malaisé, V., Zweigenbaum, P., and Bachimont, B. (2004). Detecting semantic relations between terms in definitions. In S. Ananadiou and P. Zweigenbaum, editors, *COLING 2004 CompuTerm 2004: 3rd International Workshop on Computational Terminology*, pages 55–62, Geneva, Switzerland.
- Miliaraki, S. and Androutsopoulos, I. (2004). Learning to identify single-snippet answers to definition questions. In *Proceedings of COLING 2004*, pages 1360–1366, Geneva, Switzerland.
- Nielsen, R. D. and Pradhan, S. (2004). Mixing weak learners in semantic parsing. In Lin and Wu (2004), pages 80–87.
- Pearson, J. (1996). The expression of definitions in specialised texts: a corpus-based analysis. In M. Gellerstam, J. Järborg, S. G. Malmgren, K. Norén, L. Rogström, and C. Pappmehl, editors, *Proceedings of the Seventh Euralex International Congress*, pages 817–824, Göteborg.
- Przepiórkowski, A., Degórski, Ł., and Wójtowicz, B. (2007a). On the evaluation of Polish definition extraction grammars. In Z. Vetulani, editor, *Proceedings of the 3rd Language & Technology Conference*, pages 473–477, Poznań, Poland.
- Przepiórkowski, A., Degórski, Ł., Spousta, M., Simov, K., Osenova, P., Lemnitzer, L., Kuboň, V., and Wójtowicz, B. (2007b). Towards the automatic extraction of definitions in Slavic. In J. Piskorski, B. Pouliquen, R. Steinberger, and H. Tanev, editors, *Proceedings of the Workshop on Balto-Slavonic Natural Language Processing at ACL 2007*, pages 43–50, Prague.
- Storrer, A. and Wellinghoff, S. (2006). Automated detection and annotation of term definitions in German text corpora. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation, LREC 2006*, Genoa. ELRA.
- Walter, S. and Pinkal, M. (2006). Automatic extraction of definitions from German court decisions. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics*, pages 20–28, Sydney, Australia.
- Xu, P. and Jelinek, F. (2004). Random forests in language modeling. In Lin and Wu (2004), pages 325–332.
- Xu, P. and Jelinek, F. (2007). Random forests and the data sparseness problem in language modeling. *Computer Speech and Language*, **21**(1), 105–152.