

Comparing learning approaches to coreference resolution. There is more to it than ‘bias’

Véronique Hoste
Walter Daelemans

VERONIQUE.HOSTE@UA.AC.BE

WALTER.DAELEMANS@UA.AC.BE

CNTS-Language Technology Group, University of Antwerp, Universiteitsplein 1, 2610 Wilrijk, Belgium

Abstract

On the basis of results on three coreference resolution data sets we show that when following current practice in comparing learning methods, we cannot reliably conclude much about their suitability for a given task. In an empirical study of the behavior of representatives of two machine learning paradigms, viz. lazy learning and rule induction on the task of coreference resolution we show that the initial differences between learning techniques are easily overruled when taking into account factors such as feature selection, algorithm parameter optimization, sample selection and their interaction. We propose genetic algorithms as an elegant method to overcome this costly optimization.

1. Introduction

A central question in machine learning research, and more specifically in machine learning of language research is to determine which are the learning algorithms best suited for a given task. Given the ‘no free lunch’ (Wolpert & Macready, 1995) theorem, this suitability has to be determined experimentally, most often in comparative experiments. In most comparative machine learning experiments, two or more algorithms are compared for a fixed sample selection, feature selec-

tion, feature representation, and (default) algorithm parameter setting over a number of trials (cross-validation), and if the measured differences are statistically significant, conclusions are drawn about which algorithm is better suited and why (mostly in terms of algorithm bias). Sometimes different sample sizes are used to provide a learning curve, and sometimes parameters of (some of the) algorithms are optimized on training data, but this is exceptional more than common practice. Many empirical findings, though illustrative, are observations on experiments in which one or two variables are alternated, but in which no overall optimization is undertaken (Mooney (1996), Lee and Ng (2002) and others).

In this paper, we show that there is a high risk that other areas in the experimental space may lead to radically different results and conclusions. We propose genetic search as an elegant method to overcome the computationally expensive optimization. Applied to the specific test case of coreference resolution, we experiment with two machine learning methods and discuss some methodological issues involved in running a comparative machine learning (of language) experiment. We will empirically show that changing any of the architectural variables (such as algorithm parameters, information sources, sample selection) can have great effects on the performance of a learning method.

The remainder of this paper is organized as follows. Section 2 introduces and motivates our test case task, coreference resolution and gives a short overview of the data sets and features. Section 3 presents the two machine learning packages which

Appearing in *Proceedings of the ICML-2005 Workshop on Meta-learning*, Bonn, Germany, 2005. Copyright 2005 by the author(s)/owner(s).

we used in our experiments and discusses different factors which can influence a machine learning (of language) experiment. In Section 4, we continue with a discussion of the feature selection, the parameter optimization and sample selection experiments. Section 5 reports on the use of a genetic algorithm for joint optimization. We conclude with some general observations in Section 6.

2. The task of coreference resolution

Coreference can be considered as the act of using a referring expression to point to some discourse entity. Written and spoken texts contain a large number of coreferential relations and a good text understanding largely depends on the correct resolution of these relations. In the following sentences, for example, it is the task of coreference resolution to link “they” to “The kidnappers”.

In 1983 Alfred Heineken and his driver were kidnapped. **The kidnappers** asked a ransom of 43 million guilders. A modest sum, **they** thought.

Machine learning approaches, especially supervised ones, have become increasingly popular for this problem: the C4.5 decision tree learner (Quinlan, 1993) as used by McCarthy (1996) and Soon et al. (2001), the RIPPER rule learner (Cohen, 1995) as in Ng and Cardie (2002) or a memory-based learner (Daelemans et al., 2002) as in Hoste (2005).

For the experiments, we selected all noun phrases in the English MUC-6 (MUC-6, 1995) and MUC-7 (MUC-7, 1998) corpora and the Dutch KNACK-2002 corpus (Hoste, 2005). In order to detect these noun phrases, we performed tokenization, part-of-speech tagging and NP chunking. On the basis of these preprocessed texts, we selected positive and negative instances for the training data. Positive instances were made by combining each anaphor with each preceding element in the coreference chain (a set of noun phrases referring to the same discourse entity). The negative instances were built by combining each anaphor with each preceding NP which was not part of any coreference chain and by combining each anaphor

with each preceding NP which was part of another coreference chain. This resulted in a highly skewed data set. For example, out of the 171,081 training instances in the MUC-6 data merely 6.6% were positive ones. Besides merging all NPs into one single train and test set, we also built 3 smaller datasets, each specialized in one NP type (pronouns, proper nouns, common nouns).

For our coreference resolution system, we used a combination of positional features (features indicating the number of sentences/NPs between the anaphor and its possible antecedent), morphological and lexical features (such as features which indicate whether a given anaphor, its candidate antecedent or both are pronouns, proper nouns, demonstrative or definite NPs), syntactic features which inform on the syntactic function of the anaphor and its candidate antecedent and check for syntactic parallelism, string-matching features which look for complete and partial matches and finally several semantic features. For a detailed overview of the features, we refer to Hoste (2005).

The validation experiments were performed using ten-fold cross-validation on the available training data. In order to have an idea of the performance on the minority class, we evaluated the results of our experiments in terms of precision, recall and F_β . Coreference resolution was selected as test case task for the experiments, since it implies a typical language learning task with many exceptional and low-frequency cases. Furthermore, as discussed earlier coreference resolution data sets are also highly skewed and consist of instances with some informative features and many uninformative ones (see for example Soon et al. (2001)).

3. A lazy and an eager learner

We experimented with two machine learning techniques on the task of coreference resolution: the lazy learning implementation TIMBL (Daelemans et al., 2002) and the eager rule induction method RIPPER (Cohen, 1995). The learning biases of these two approaches provide extremes in the *eagerness* dimension in ML (the degree in which a learning algorithm abstracts from the training data in forming a hypothesis). The motivation for

the choice of a lazy or an eager learning bias may be understandability of learned models or abstraction from noise (eager) or the possibility of learning from low-frequency or untypical data points (lazy) to name but a few. But comparing two or more algorithms on a specific task is complex. In Figure 1, a characterization of this complexity is given.

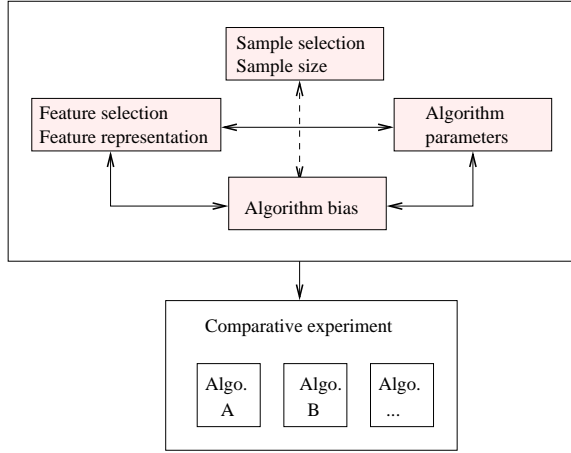


Figure 1. Graphical representation of the aspects influencing a (comparative) machine learning experiment. The filled lines and the dashed line represent the experiments reported in this paper. The dashed line also refers to previous research on sample size of Banko and Brill (2001).

Apart from the algorithm bias, many other factors potentially play a role in the outcome of a machine learning experiment. One of these factors is **the data set**, viz. the sample selection and its size (Banko & Brill, 2001). Also the selection of high-quality training instances has an important effect on predictive accuracy. Furthermore, class imbalances in the selected data set can also affect classification results. Another influential factor are the **information sources** used: the **features selected** for prediction, and **their representation** (e.g. binary, numeric or nominal). The presence of irrelevant features can considerably slow the rate of learning and have a negative effect on classification results. Furthermore, most learning algorithms have a number of **algorithm parameters** which can be tuned. These factors also interact: a feature selection which is optimal with default parameter settings is not necessarily optimal when changing the algorithm parameters. The optimal algorithm parameters for a

skewed data set will not necessarily be optimal when changing the class distribution in the data.

In the following section, we will show that performance differences due to algorithm parameter optimization, feature selection, and sample selection can easily overwhelm the performance differences reported between algorithms in comparative experiments. Due to the large number of experiments performed, we will mainly discuss the observed tendencies which hold for all data sets.

4. The effect of optimization

4.1. Searching the feature space

Although the search for disambiguating features is central in the machine learning research for coreference resolution and for NLP tasks in general, there is no general practice to also consider the complex interaction between all these information sources. For our experiments, we used two automated techniques for the selection of the relevant features, viz. backward elimination (John et al., 1994) and bidirectional hill-climbing (Caruana & Freitag, 1994). Table 1 gives the results of these experiments for TIMBL and RIPPER on the MUC-7 data sets. It shows that (i) the algorithm-comparing differences can be overruled by the algorithm-internal performance differences and that (ii) especially TIMBL can benefit from feature selection which is mainly due to the embedded feature selection in the construction of the rules in RIPPER and the fact that TIMBL does not take into account dependencies between features.

With respect to the selected features, no general conclusions could be drawn. Per data set and per selection procedure, a different feature set is selected by each learner, which implies that the optimal feature selection has to be determined experimentally for each single data set.

4.2. The effect of parameter optimization

Another factor which can have great effects on classifier performance is the choice of algorithm parameter settings. Although both algorithms provide sensible default settings, it is by no means certain that they are the optimal settings for our

Table 1. Results of TIMBL and RIPPER on the MUC-7 data sets after backward selection and bidirectional hill-climbing.

		TIMBL			RIPPER		
		Prec.	Rec.	$F_{\beta=1}$	Prec.	Rec.	$F_{\beta=1}$
All	default	51.57	46.09	48.68	77.51	36.21	49.36
	backward	73.98	41.26	52.98	77.49	40.34	53.06
	bi.hill.	75.39	42.08	54.01	77.99	40.52	53.33
Pronouns	default	42.31	36.60	39.25	59.50	22.70	32.86
	backward	46.86	40.89	43.67	59.34	26.43	36.57
	bi.hill.	61.55	25.51	36.07	60.74	30.94	41.00
Proper nouns	default	62.36	56.87	59.49	84.58	52.56	64.83
	backward	73.92	55.54	63.43	87.08	55.22	67.59
	bi.hill.	85.18	54.88	66.75	88.20	51.72	65.21
Common Nouns	default	43.06	39.17	41.03	74.56	36.76	49.24
	backward	52.02	39.28	44.76	76.05	40.41	52.78
	bi.hill.	78.83	38.72	51.93	76.77	39.70	52.33

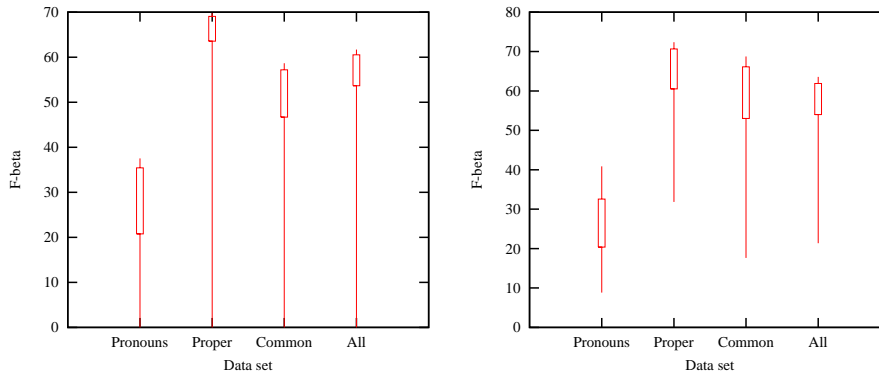


Figure 2. Results of TIMBL (left) and RIPPER (right) over all parameter settings for all MUC-6 data sets. The graphs show the difference between the performance obtained with the best and worst parameter settings per data set. The boxes in the graphs represent averages and deviations.

task. Therefore, we exhaustively varied the algorithm parameter settings for each classifier. For TIMBL, the following parameters were varied: the similarity metric (overlap or modified value difference metric (MVDM)), feature weighting (no weighting, information gain weighting, gain ratio weighting, ...), neighbor weighing (majority voting and different distance weighting schemes), and the k parameter which controls the number of nearest neighbors. For RIPPER, there was an optimization of the class ordering parameter (+freq, -freq or mdl), the two-valued negative tests parameter (-ln or nothing), the hypothesis simplification parameter (0.5, 1 or 1.5), the example coverage parameter (1, 2, 3 or 4), the parameter expressing the number of optimization passes (0, 1 or 2) and the loss ratio parameter (0.5, 1, 1.5). Figure 2 displays the $F_{\beta=1}$ results of all algo-

rithm parameter optimization experiments for the MUC-6 data sets. Per algorithm and per data set, the best and worst scores are displayed, as well as the averages and deviations.

The long vertical lines in Figure 2 reveal a lot of variation in the $F_{\beta=1}$ results when varying the algorithm parameters, although the boxes which are mostly located in the upper area indicate that the badly performing parameter combinations are in the minority. A good parameter combination is crucial. In the MUC-6 common nouns data set, for example, RIPPER yields an $F_{\beta=1}$ score of 17.65% for the combination of the ‘-freq’ ordering method and the ‘0.5’ loss ratio value. Combining this below zero loss ratio value with the ‘+freq’ or ‘mdl’ ordering methods, however, leads to top $F_{\beta=1}$ scores on the validation material. Similar observations can be made for TIMBL.

Overall, we observed that parameter optimization leads to large performance increases for both learners. Furthermore, we observed that parameters cannot be generalized. The optimal settings merely reveal some tendencies.

4.3. The effect of sample selection

Coreference resolution data sets reveal large class imbalances: only a small part of the possible relations between noun phrases is coreferential. In the KNACK-2002 cross-validation data, for example, merely 6.3% of the instances is classified as positive. Learning performance can be hindered when learning from these data sets where the minority class is underrepresented. A central question in the discussion on data sets with an imbalanced class distribution is in what proportion the classes should be represented in the training data. In the machine learning literature, there have been several proposals (see Japkowicz and Stephen (2002)) for adjusting the number of majority class and minority class examples. Methods include resizing training data sets or sampling, adjusting misclassification costs, learning from the minority class, adjusting the weights of the examples, etc.

In order to investigate the effect of class distribution on classifier performance, we compared the performance of the classifiers on a variety of class distributions. We investigated the effect of random down-sampling and down-sampling of the true negatives for both TIMBL and RIPPER. This was done by gradually downsizing the number of negatives instances in slices of 10% until there was an equal number of positive and negative training instances. These experiments reveal the same tendencies for the different data sets. As exemplified in Figure 3, we can observe that TIMBL and RIPPER behave differently. RIPPER is more sensitive to the skewedness of the classes and down-sampling is beneficial for the RIPPER results. Furthermore, down-sampling only starts being harmful at a high down-sampling level. TIMBL has shown this tendency only on the ‘‘Pronouns’’ data set. But no down-sampling level leads to the best performance over all data sets.

As an illustration of RIPPERS sensitivity to

skewedness, we also varied the loss ratio parameter in RIPPER, which allows the user to specify the relative cost of the false positives and false negatives. In its default version, RIPPER uses a loss ratio of 1, which indicates that the two errors have equal costs. In order to improve on recall, we varied the loss ratio in RIPPER from 1 (default) to 0.05. As shown in Table 2 for the different KNACK-2002 data, a change of the loss ratio parameter leads to a large performance increase over the default settings.

Table 2. Default $F_{\beta=1}$ results for the KNACK-2002 data, in comparison with the highest and lowest scores after change of the loss ratio parameter (between brackets).

	default	high	low
All	46.49	60.33 (0.2)	46.49
Pronouns	50.57	63.49 (0.4)	50.57
Proper nouns	60.21	63.69 (0.06)	58.61
Common nouns	36.52	42.68 (0.07)	36.52

With respect to the specific loss ratio values, we conclude that no particular value leads to the best performance over all data sets. This confirms our findings in the down-sampling, parameter optimization and feature selection experiments, which also revealed that the optimal parameters and features for a given task have to be determined experimentally per data set.

5. Genetic algorithms for joint optimization

In a final optimization step we explored the interaction between the previously mentioned factors. Joint feature selection, sample selection and parameter optimization is essentially an optimization problem which involves searching the space of all possible feature subsets, sample subsets and parameter settings to identify the combination that is optimal or near-optimal. Given the combinatorially explosive character of this type of joint optimization, we have chosen for genetic algorithms (GA, e.g. Goldberg (1989) and Mitchell (1996)) as a computationally feasible way to achieve this. One of the advantages of genetic algorithms in contrast to local search methods such as hill-climbing, gradient based and simu-

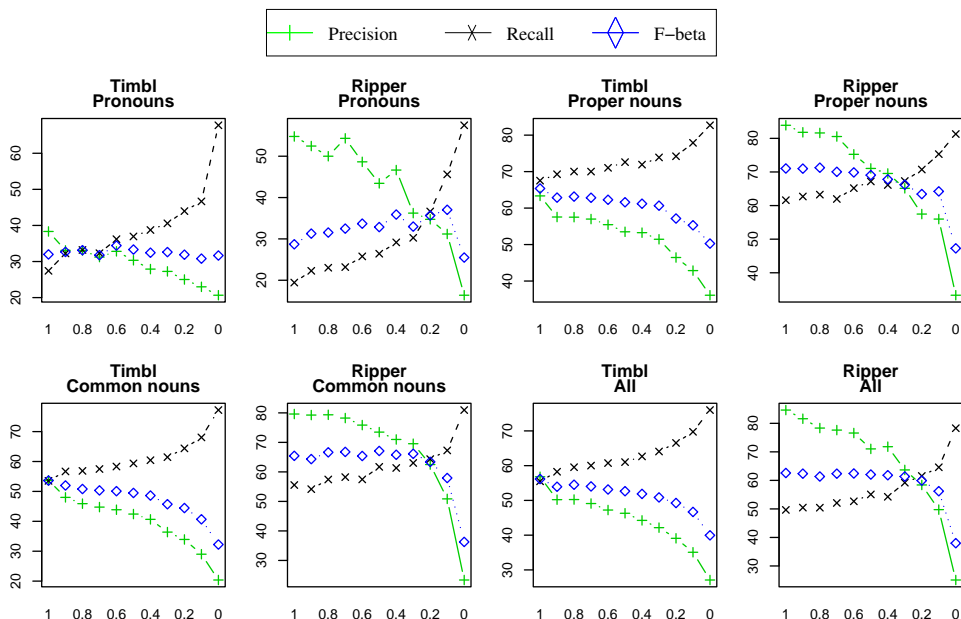
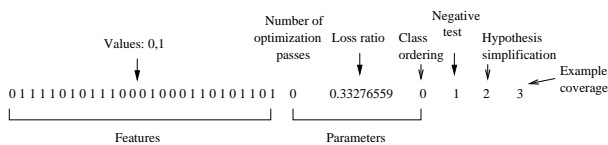


Figure 3. Cross-validation results in terms of precision, recall and $F_{\beta=1}$ after application of TIMBL and RIPPER on the MUC-6 data with a randomly down-sampled majority class. The test partitions keep their initial class distribution.

lated annealing methods is that they explore different areas of the search space in parallel, which might be a reasonable strategy for problems with a large number of parameters and features. GAs contain at any time a population of candidate solutions to the optimization problem to be solved. For the experiments, we used a generational genetic algorithm with the evaluations distributed over a cluster of computers using the Sun Grid Engine queuing system. The following GA parameters were used and kept constant: maximal number of generations=30, population size=10, uniform crossover (crossover rate=0.9), tournament selection (selection size=2), discrete (mutation rate=0.2) and Gaussian mutation (k and loss ratio). $F_{\beta=1}$ was used as fitness function. We are aware that the optimization problem we are trying to solve with a genetic algorithm also applies to the GA itself.



In the experiments, the individuals were represented as bit strings. Each individual contains

particular values for all algorithm settings and for the selection of the features. For example, for RIPPER, the features are encoded as binary alleles. At the end of the chromosome, the different algorithm parameters are represented. Through the variation of the loss ratio parameter, which controls the relative weight of precision versus recall, a down-sampling effect can be obtained.

The GA optimization experiments confirm the tendencies observed in the optimization experiments in the previous section. As exemplified for the KNACK-2002 data in Figure 4, we could observe for all data sets that the performance differences inside one single learning method can be much larger than the method-comparing performance differences. In their default representation, for example, TIMBL and RIPPER yield a 46.8% and a 46.5% $F_{\beta=1}$ score, respectively. Optimization leads to a large performance improvement for both learners and to a reversed supremacy: 55.7% for TIMBL and 61.7% for RIPPER. In conclusion, we can state that we cannot draw conclusions of one classifier being better on a particular task than another classifier, when only taking into account default settings or limited optimization.

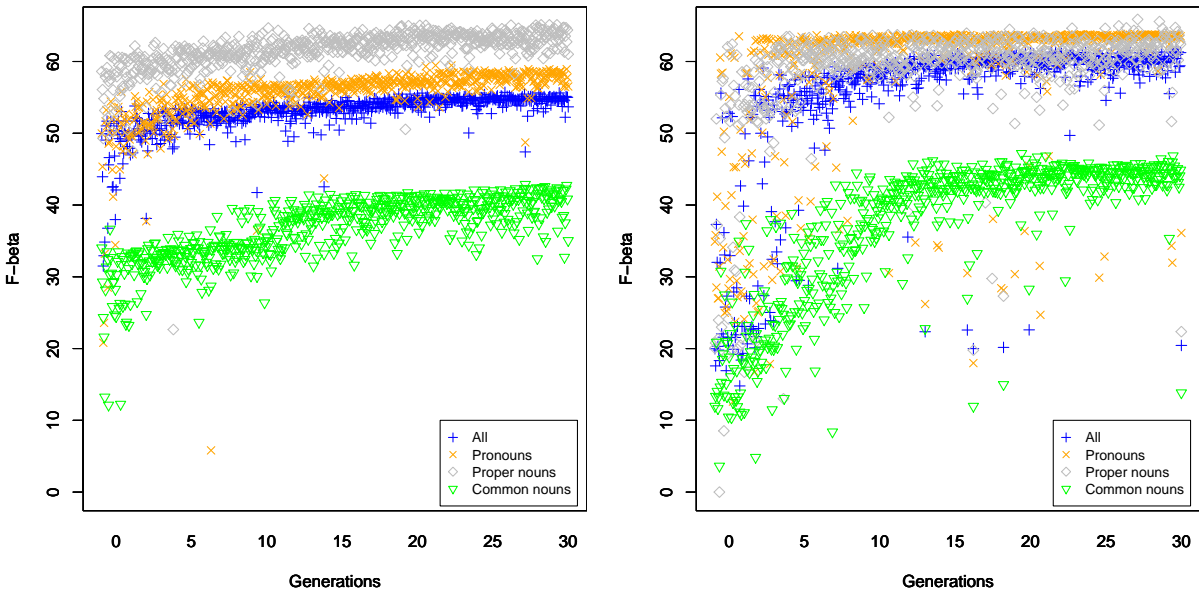


Figure 4. Results of TIMBL (left) and RIPPER (right) for the KNACK-2002 data sets during GA optimization.

These observations, however, are not limited to the task of coreference resolution. In earlier work (Hoste et al., 2002; Daelemans et al., 2003), we came to similar conclusions for the task of word sense disambiguation, the prediction of diminutive suffixes and part-of-speech tagging. Furthermore, this effect of optimization is not limited to natural language processing datasets. We performed experiments on 5 UCI benchmark datasets¹ and we came to similar conclusions as on the NLP data sets. These effects explain why in the machine learning of natural language literature, so many results and interpretations about the superiority of one algorithm over the other are contradictory. We show that there is a high risk that other areas in the experimental space may lead to radically different results and conclusions.

¹<http://www.ics/uci/edu/~mlearn/MLRepository.html>. The experiments were performed on “database for fitting contact lenses” (24 instances), “contraceptive method choice” (1473 instances), “breast-cancer-wisconsin” (699 instances), “car evaluation data base” (1728 instances) and “postoperative patient data” (90 instances).

6. Concluding remarks

In this paper, we have investigated the behavior of two classification-based learning approaches when learning coreference resolution. We showed that many factors can affect the success of a classifier, such as the specific ‘bias’ from the classifier, the choice of algorithm parameters, the selection of information sources, the sample selection and the interaction between these factors. We also showed that optimization can lead to radically different results, causing much larger classifier-internal variations than classifier-comparing variations. These results call into question the usefulness of the numerous classifier comparison studies in the literature. On the other hand, significant performance increases can be obtained this way. We conclude that in general, the more effort is put in optimization, through feature selection, parameter optimization, sample selection and their joint optimization, the more reliable the results and the comparison will be.

References

- Banko, M., & Brill, E. (2001). Scaling to very very large corpora for natural language disambiguation. *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics* (pp. 26–33).
- Caruana, R., & Freitag, D. (1994). Greedy attribute selection. *Proceedings of the International Conference on Machine Learning* (pp. 28–36).
- Cohen, W. W. (1995). Fast effective rule induction. *Proceedings of the 12th International Conference on Machine Learning* (pp. 115–123).
- Daelemans, W., Hoste, V., De Meulder, F., & Naudts, B. (2003). Combined optimization of feature selection and algorithm parameter interaction in machine learning of language. *Proceedings of the 14th European Conference on Machine Learning* (pp. 84–95).
- Daelemans, W., Zavrel, J., van der Sloot, K., & van den Bosch, A. (2002). *Timbl: Tilburg memory-based learner, version 4.3, reference guide* (Technical Report ILK Technical Report - ILK 02-10). Tilburg University.
- Goldberg, D. (1989). *Genetic algorithms in search, optimization and machine learning*. Addison Wesley.
- Hoste, V. (2005). *Optimization issues in machine learning of coreference resolution*. Doctoral dissertation, Antwerp University.
- Hoste, V., Hendrickx, I., Daelemans, W., & van den Bosch, A. (2002). Parameter optimization for machine-learning of word sense disambiguation. *Natural Language Engineering, Special Issue on Word Sense Disambiguation Systems*, 8, 311–325.
- Japkowicz, N., & Stephen, S. (2002). The class imbalance problem: A systematic study. *Intelligent Data Analysis Journal*, 6, 429–450.
- John, G., Kohavi, R., & Pfleger, K. (1994). Irrelevant features and the subset selection problem. *International Conference on Machine Learning* (pp. 121–129).
- Lee, Y., & Ng, H. (2002). An empirical evaluation of knowledge sources and learning algorithms for word sense disambiguation. *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing* (pp. 41–48).
- McCarthy, J. (1996). *A trainable approach to coreference resolution for information extraction*. Doctoral dissertation, Department of Computer Science, University of Massachusetts, Amherst MA.
- Mitchell, M. (1996). *An introduction to genetic algorithms*. MIT Press.
- Mooney, R. (1996). Comparative experiments on disambiguating word senses: An illustration of the role of bias in machine learning. In E. Brill and K. Church (Eds.), *Proceedings of the conference on empirical methods in natural language processing*, 82–91.
- MUC-6 (1995). Coreference task definition. version 2.3. *Proceedings of the Sixth Message Understanding Conference* (pp. 335–344).
- MUC-7 (1998). Muc-7 coreference task definition. version 3.0. *Proceedings of the Seventh Message Understanding Conference*.
- Ng, V., & Cardie, C. (2002). Combining sample selection and error-driven pruning for machine learning of coreference rules. *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing* (pp. 55–62).
- Quinlan, J. (1993). *C4.5: Programs for machine learning*. Morgan Kaufmann, San Mateo, CA.
- Soon, W., Ng, H., & Lim, D. (2001). A machine learning approach to coreference resolution of noun phrases. *Computational Linguistics*, 27, 521–544.
- Wolpert, D., & Macready, W. (1995). *No free lunch theorems for search* (Technical Report SFI-TR-95-02-010). Santa Fe Institute, Santa Fe, NM.