

RESEARCH STATEMENT

Roi Reichart (roiri@cs.huji.ac.il, www.cs.huji.ac.il/~roiri)

Natural Language processing (NLP) is a field that combines linguistics, cognitive science, statistical machine learning and other computer science areas in order to compile intelligent computer systems that can understand human languages. NLP has various applications, among which are machine translation, question answering and search engines.

The field of NLP has, in the past two decades, come to simultaneously rely on and challenge the field of machine learning. Statistical methods now dominate NLP, and have moved the field forward substantially, opening up new possibilities for the exploitation of data in developing NLP components and applications.

Many state of the art natural language algorithms are based on supervised learning techniques. In this type of learning, a corpus consisting of texts annotated by human experts is compiled and used to train a learning algorithm. While supervised learning has made substantial contribution to NLP, it faces some significant challenges.

Many fundamental NLP tasks, such as syntactic parsing, part-of-speech (POS) tagging and machine translation, involve structured prediction and sequential labeling. For such kind of tasks, compiling annotated corpora is costly and error prone due to the complex nature of annotation. I refer to this challenge as the annotation bottleneck.

A closely related challenge is that of *domain adaptation*. Supervised algorithms usually perform well when the training data and the data for which they should provide predictions (the test data) are drawn from similar domains. When an algorithm trained with data from one domain is to provide predictions for data taken from a substantially different domain, its performance markedly degrades. Creating corpora for every test domain is not feasible due to the aforementioned annotation bottleneck.

Supervised natural language learning is also challenged by methodological problems. Annotation schemes for tasks such as syntactic parsing are often based on arbitrary decisions. Such schemes often provide a detailed description of certain structures while addressing others only briefly. Many applications would benefit from different annotation decisions.

My research focuses on developing machine learning techniques that deal with these challenges. Its main theme is utilizing the plentiful amounts of raw text available nowadays, for creating state of the art algorithms that use little to no manually annotated data. Cases where little amounts of manually annotated text is utilized are referred to as *semi-supervised learning* [1]; while cases where only raw text is used are known as *unsupervised learning* [2, 3]. I explore semi-supervised and unsupervised techniques for structured prediction and sequence labeling NLP tasks such as syntactic parsing and POS tagging.

Combining Labeled and Unlabeled Data: Semi-Supervised Learning

My research of semi-supervised learning has followed three main directions:

1. Utilizing raw text to enhance the performance of syntactic parsers when only small amounts of labeled data are available [4]. When parsers are used by applications in the highly variable web, this is the rule rather than the exception.
2. Selecting high quality text samples for statistical parsers training, whose annotation requires minimal human effort. Active learning [5] has traditionally been used for this task [6]. My

research focuses on finding both reliable measures for the human annotation efforts and developing algorithms for selecting the samples [7, 8].

3. Developing reliable confidence scores for the output of statistical parsers [9, 10]. The output of statistical parsers is often used by NLP applications such as question answering and machine translation. A confidence score can provide these applications with an indication about the effect that a given syntactic analysis will have on their performance. In addition, confidence scores play an important role in many semi-supervised learning protocols [11].

For the first direction, various semi-supervised techniques have been explored throughout the years. Among these are active-learning, co-training and ensemble methods. Our work was the first to show that self-training, a technique where the parser is trained on its own output, can substantially enhance parser performance in the small annotated training data case. The self-training protocol we proposed, currently achieves the best results for this scenario.

For the second direction, we have addressed two scenarios that have not been explored in the literature before. One case is where a corpus is to be built for multiple linguistics annotation schemes. In our work we developed algorithms for simultaneous annotation with syntactic trees and name entities.

The other case is that of building a syntactically annotated corpus where the cognitive efforts of the human annotator are to be minimized while the corpus quality as training material for statistical parsers is to be maximized. In this work we proposed measures for the human cognitive efforts in syntactically analyzing a sentence and developed the algorithms accordingly. Our algorithm has been shown to have better performance on the task when compared to previously proposed algorithms. The paper that summarizes this research received the **best paper award** in CoNLL 2009.

For the third direction, we developed state of the art confidence measures for both supervised and unsupervised parsers. In the supervised case the measure is based on an ensemble technique while in the unsupervised case it utilizes a cluster-based tree representation. In both cases, the proposed score has been shown to outperform other existing scores.

In addition to its value for applications, the score can indicate if the tree induced by the parser can be used as training data for supervised and unsupervised parsers. This potential application of quality scores is especially important in the unsupervised case, since high quality trees produced by unsupervised algorithms can provide training material for supervised parsers with no human efforts. Our recent results indicate that the unsupervised quality score can be used in a way that substantially improves unsupervised parser performance.

Language Acquisition from Plain Text: Unsupervised Learning

In the unsupervised direction, my research has focused on sequential clustering tasks such as POS tagging and the induction of labeled parse trees. We have explored both algorithmic approaches and the unique evaluation challenges of the field.

Labels induced by unsupervised algorithms receive arbitrary names. Evaluating their quality against a gold standard therefore requires a correspondence scheme between the induced and gold labels. We have explored mapping-based and information-theoretic-based approaches for this challenge and analyzed their properties [12, 13].

For the computation of mapping-based evaluation schemes we explored graph-based optimization methods such as the Kuhn-Munkres algorithm. For the information theoretic evaluation

schemes [14] we developed measures based on conditional and unconditional entropies and analyzed their role in reflecting the clustering quality. We used these measures to analyze existing algorithmic approaches to the aforementioned sequential clustering problems, and to suggest novel techniques.

A property shared by most algorithmic approaches to the aforementioned tasks is that the optimal clustering is defined to be the one that optimizes a non-convex function of the data. Optimization algorithms for such functions are well-known to converge to local maxima which yield clustering solutions of variable quality. Finding a final, high quality solution is a difficult problem that has hardly been addressed in the NLP literature. In a recent work we developed a method for identifying the high quality runs of unsupervised POS induction algorithms [15].

We proposed various algorithmic techniques that deal with this problem. For the induction of labeled parse trees we developed clustering techniques based on a Bayesian minimum description length (MDL) principle [17] and on ensembles of Dirichlet Process mixture model experts. For POS induction we developed both a deterministic algorithm based on linguistics distributional and morphological considerations and a method that imposes existing stochastic methods to converge to a high quality solution [16, 15]. All these methods output a high quality, unique final clustering (best reported results for the addressed problems).

My research on semi-supervised and unsupervised learning follows a few characteristics guidelines. It explores methods for the application of unlabeled text in a statistically robust manner and aims at developing evaluation measures that reflect the various statistical and cognitive dimensions of natural language acquisition.

Research Agenda: Directions for Future Work

The research I performed during my PhD is obviously just a first step in the journey towards building computer algorithms that are able to acquire human languages from plain text without the guidance of a human expert.

On the algorithmic level, efficient algorithms that are based on sophisticated statistical models and can still be applied to plentiful amounts of raw text are to be developed. Moreover, our understating of what are human languages is to be enhanced. One step in this direction is using unsupervised methods such as those developed in my research for the design of data-driven tag sets for NLP tasks, hopefully reflecting the nature of the data better than manually designed tag sets.

In a more general perspective, NLP research is strongly related to the development of advanced learning techniques. As can be inferred from the research described above, the relevant fields are domain adaptation, unsupervised methods and structured prediction. I intend to devote efforts to these machine learning fields in the future.

In my research I also explored the application of syntactic acquisition algorithms to semantic acquisition tasks such as ontology induction [18], verb argument identification [19], learning the syntactic/semantic frame of a verb [20] and understanding the role of commas in sentences [21]. This research can potentially cast light on the role of syntax in defining the semantic meaning of sentences. Moreover, it can yield novel measures for the output quality of syntactic acquisition algorithms based on the value of their output for semantic acquisition. Deeper understanding of the connection between syntax and semantics is an important topic for future research.

During my PhD research I have published 14 papers in the major international NLP conferences (ACL, COLING and CoNLL). One of these papers received the **best paper award** in CoNLL 2009; five more works are under review for COLING 2010.

Ten of these nineteen papers are the fruit of cooperation with at least one co-author (that is, in addition to my PhD supervisor, Prof. Ari Rappoport, who is a co-author of all my papers), and the total number of co-authors with which I cooperated is ten, 5 of which work outside Israel (3 in US, 2 in German). Cooperating with other people is a fundamental aspect of my research style. It enhances my work and enables it to face more difficult challenges and address more complex issues. Of course, I find it a much enjoyable way of work.

References

- [1] Steven Abney, “Semisupervised Learning for Computational Linguistics.”, *Chapman and Hall*, 2007.
- [2] Dan Klein, “The unsupervised learning of natural language structure.”, *PhD Thesis, Stanford University*, 2005.
- [3] Alexander Clark, “Unsupervised Language Acquisition: Theory and Practice.”, *PhD Thesis, University of Sussex*, 2001.
- [4] **Roi Reichart** and Ari Rappoport, “Self-Training for Enhancement and Domain Adaptation of Statistical Parsers Trained on Small Datasets.”, *ACL ’07*.
- [5] David Cohn, Les Atlas and Richard Ladner. “Improving generalization with active learning.”, *Machine Learning*, 15(2):201–221.
- [6] Rebecca Hwa, “Sample selection for statistical parsing.”, *Computational Linguistics*, 30(3):253–276.
- [7] **Roi Reichart**, Katrin Tomanek, Udo Hahn and Ari Rappoport, “Multi-Task Active Learning for Linguistic Annotations.”, *ACL ’08*.
- [8] **Roi Reichart** and Ari Rappoport, “Sample Selection for Statistical Parsers: Cognitively Driven Algorithms and Evaluation Measures.”, *CoNLL ’09*.
- [9] **Roi Reichart** and Ari Rappoport, “An Ensemble Method for Selection of High Quality Parses.”, *ACL ’07*.
- [10] **Roi Reichart** and Ari Rappoport, “Automatic Selection of High Quality Parses Created By a Fully Unsupervised Parser.”, *CoNLL ’09*.
- [11] Avrim Blum and Tom Mitchell, “Combining labeled and unlabeled data with co-training.”, *COLT ’98*.
- [12] **Roi Reichart** and Ari Rappoport, “The NVI Clustering Evaluation Measure.”, *CoNLL ’09*.
- [13] **Roi Reichart**, Omri Abend and Ari Rappoport, “Type Level Clustering Evaluation: New Measures and a POS Induction Case Study.”, *CoNLL ’10*.
- [14] Marina Meila, “Comparing clustering – an information based distance.”, *Journal of Multivariate Analysis*, 98:873–895.
- [15] **Roi Reichart**, Raanan Fattal and Ari Rappoport, “Improved Unsupervised POS Induction Using Intrinsic Clustering Quality and a Zipfian Constraint.”, *CoNLL ’10*.
- [16] Omri Abend, **Roi Reichart** and Ari Rappoport, “Improved Unsupervised POS Induction through Prototype Discovery.”, *ACL ’10*.
- [17] **Roi Reichart** and Ari Rappoport, “Unsupervised Induction of Labeled Parse Trees by Clustering with Syntactic Features.”, *COLING ’08*.
- [18] Dmitry Davidov, **Roi Reichart** and Ari Rappoport, “Superior and Efficient Fully Unsupervised Pattern-based Concept Acquisition Using an Unsupervised Parser.”, *CoNLL ’09*.
- [19] Omri Abend, **Roi Reichart** and Ari Rappoport, “Unsupervised Argument Identification for Semantic Role Labeling.”, *ACL ’09*.
- [20] Omri Abend, **Roi Reichart** and Ari Rappoport, “A Supervised Algorithm for Verb Disambiguation into Verb-Net Classes.”, *COLING ’08*.
- [21] Vivek Srikumar, **Roi Reichart**, Mark Sammons, Ari Rappoport and Dan Roth, “Extraction of Entailed Semantic Relations Through Syntax-based Comma Resolution.”, *ACL ’08*.