

Towards Automatic Extraction of Definitions

Claudia Borg¹, Michael Rosner¹, and Gordon J. Pace²

¹ Department of Artificial Intelligence, University of Malta

² Department of Computer Science, University of Malta
{claudia.borg|mike.rosner|gordon.pace}@um.edu.mt

Abstract. Definition extraction can be useful for the creation of glossaries and in question answering systems. It is a tedious task to extract such sentences manually, and thus an automatic system is desirable. In this work we review various attempts at rule-based approaches reported in the literature and discuss their results. We also propose a novel experiment involving the use of genetic programming and genetic algorithms, aimed at assisting the discovery of grammar rules which can be used for the task of definition extraction.

1 Introduction

A definition is a term together with a description of its meaning or the concept it refers to. Definitions are helpful because they facilitate the understanding of new terms. The extraction of definitions from text can be useful in various scenarios, including the automatic creation of glossaries for the building of dictionaries and in question answering systems. In this work, we will focus on the use of definition extraction in eLearning, where definitions can help learners conceptualise new terms and help towards the understanding of new concepts encountered in learning material.

eLearning is the process of acquiring knowledge through electronic aids, by providing access to materials that will enable them to learn a particular task. A tutor can direct the learning process through a Learning Management System (LMS), where learning material is presented to the student according to the tutor's direction. The tutor can also use the LMS to add new content, create courses, structure layout and presentation of courses and monitor student performance.

Learning material, packaged into units known as learning objects (LOs), normally contain implicit information in natural language, which would require a lot of work for the tutor to extract manually. An LMS can be enhanced by introducing tools to extract such information automatically. One such piece of information is the presence of definitions in texts. We propose a tool that will attempt to extract definitions from LOs, which an LMS could use to extract definitions (for instance, to be used to create a glossary). The tutor will then be able to refine this information rather than create it from scratch.

The task of definition extraction is a challenging one. We are trying to identify sentences that contain knowledge which could then be used by applications such

as those mentioned above. What more, we are attempting to identify sentences which *define* a term, rather than simply describe it vaguely, or compare it to other terms.

Since definitions are made up of natural language texts, we propose to use linguistic knowledge such as part-of-speech tagging and morphological analysis to support the definition extraction. Furthermore, it was noticed that there exist different syntactic forms of definitions. Hence, rather than trying to identify arbitrary definitions, we propose to look at the different definition categories separately.

We propose an experiment in Section 2 which combines genetic algorithms (GA) and genetic programming (GP) to try and discover grammar rules that could identify definitions present in LOs. Section 3 will review the work related to our task, divided into two parts. First we look at different approaches in definition extraction using rule-based techniques. Then we will look at the area of grammatical inference and the application of GAs and GPs to learn grammars. The outcome of this work is to evaluate the use of machine learning techniques (GAs and GPs) and their results in learning restricted grammars. The grammars developed through these experiments can then be applied by rule-based techniques to extract definitions. The results of the GP and the GA will be used to discover features which identify certain definitions with a high rate of accuracy, but also other features to classify less clearcut definitions using features in a combined manner.

This work is done in collaboration with an EU-funded FP6 project LT4eL³. The project is described in more detail in [BR07] and [MLS07], and aims at enhancing LMSs by using language technologies and semantic knowledge.

2 Proposed Approach

In eLearning, LOs are generally created by tutors in different formats such as HTML, PDF and other text formats. A corpus of LOs, gathered within the LT4eL project, has been converted to one standard XML format with added linguistic information. Work carried out in the project has also manually identified and annotated a set of 450 definitions from this corpus.

Given that a corpus contains both a set of definitions and a (usually larger) set of non-definitions, an attempt to learn the importance of features present in definitions is possible. A feature can be seen as a description of characteristics that can help us identify a definition. To simplify the identification process, definitions have been split into six different categories (as described in [BR07]). By learning to identify the categories separately, we reduce the size and complexity of the search space.

³ Language Technologies for Learning www.lt4e1.eu

2.1 Experiment One: Genetic Algorithm

A genetic algorithm (GA) [Hol75,Gol89] is a possible technique that can be used to learn the importance of the features that can recognise definitions. This can be done by assigning weights to each feature and allowing the algorithm to adjust the weights according to the performance. It also makes it ideal to run the GA on the separate categories of definitions identified (as described in [BR07]), so that the results can be directed to one given situation at a time.

A feature is a function which given a sentence (which includes linguistic information) returns a numeric score. An example of a feature would be that of a part-of-speech sequence that may capture a definition — returning a numeric value indicating how close the given sentence matches that particular sequence (e.g. DT→NN→VBZ→DT→NN→IN→NNS). An other example of a feature is a test whether the sentence contains the verb ‘to be’ — with the only possible values of the score now being 1 or 0, indicating the presence or otherwise of the verb. A feature is said to be an effective classifier of definitions if it gives higher scores to sentences which define a term than to other sentences.

Given a set of n basic features, f_1 to f_n , and n numeric constants, α_1 to α_n , one can produce a new feature combining these basic features in a linear fashion:

$$F(s) = \sum_{i=1}^n \alpha_i \times f_i(s)$$

The problem is now: given a fixed set of features, how can we calculate a good set of weights so as to maximise the effectiveness of the combined features as a definition classifier? We propose to use a GA to learn these weights. The values learnt would thus correspond to the relative effectiveness of the individual features as classifiers of definitions. Before starting the experiment, a predefined set of features will be adopted, which will remain static throughout the experiment. A gene will be a list of length equal to the number of predefined features of numbers. Thus, the i th gene ($1 \leq i \leq \text{population size}$) will have the following structure:

$$g_i = \langle \alpha_{i,1}, \alpha_{i,2} \dots \alpha_{i,n} \rangle$$

Note that n corresponds to the number of predefined features. The interpretation of the gene is thus a list of weights which will be applied to the output of the predefined features. For instance, $\alpha_{i,1}$ is the weight that the gene g_i would assign to the first feature. Such a gene will therefore return a score to a given sentence s as follows:

$$\text{score}_i(s) = \sum_{j=1}^n f_j(s) \times \alpha_{i,j}$$

The initial population will consist of genes that contain random weights assigned to each feature. The interpretation of the gene is a function that when applied

to a sentence gives the summation of the feature-function scores multiplied by their weights.

The fitness function will take an individual and produce a score according to its performance. The score will be calculated by applying the gene to both positive and negative examples and will be judged according to how the gene is able to separate the two sets of data from each other.

Crossover and mutation will be carried out in the traditional way of GA. Crossover will take two individuals, split them at a random position and create two new children by interchanging the parts of the parents. Mutation will take a random position in the gene and change its value. If the children perform better than the parents, they will replace them.

Once the population converges, the weights of the best individual would give the clearest separating threshold between definitions and non-definitions. This will also allow us to identify which of the features in the predefined feature set are most important.

2.2 Experiment Two: Genetic Programming

Genetic Programming (GP) [Koz92] is a technique that uses GA principles to evolve simple programs. The main difference between GAs and GPs is the representation of the population and how the operations of crossover and mutation are carried out. The members of the population are parse trees, usually of computer programs, whose fitness is determined by execution. Crossover and mutation are carried out on subtrees, ensuring that the resulting tree would still be of the correct structure.

Whereas the scope of the previous experiment was to learn which are the best performing features for the task of definition extraction from a set of predefined ones, this experiment aims at identifying new features. The choice of what type of structure we are trying to learn determines the complexity of the search space. In our application, two possible options could be regular languages (in the form of regular expressions) or context-free languages (in the form of context-free grammars), the latter having a larger search space than the former. Through observations and current work with *lxtransduce* [Tob05], regular expressions (extended with a number of constructs) should be sufficient to produce expressions that would, in most cases, correctly identify definitions.

Both basic and combined features used in the GA can serve to inject an initial population into the GP. The selection can be made based on the weights learned by the GA and translating those features into extended regular expressions. The extensions that are being considered are conjunction (possibly only at the top level due to complexity issues), negation and operators such as *contains sub-expression*. Note that some of these can already be expressed as regular expressions, however, introducing them as a new single operator helps the genetic program learn faster.

The population will evolve with the support a fitness function in order to select those individuals for mating. The fitness function can apply the extended regular expressions on the given training set and then use measurements such as precision

and recall over the captured definitions. Such measurements can indicate the performance of the individuals and will allow us to fine-tune the GP according to the focus of the experiment (where one could emphasis on a high percentage for one measurement at a time, or take an average for both). This flexibility will also allow us to have different results in the various runs of the experiment, where, for instance, in one we could try to learn over-approximations whereas in an other we can learn an under-approximation.

Crossover takes two trees and creates two new children by exchanging nodes with a similar structure. If an offspring is able to parse correctly one definition, it survives into the next generation, otherwise it is discarded. Parents would normally also be retained in the population, since we would not want to lose the good individuals (it is not obvious that their offspring would have the same capability of correctly identifying definitions).

Mutation takes an individual and selects at random one node. If that node is a leaf, it is randomly replaced by a terminal symbol. If it is an internal node, it is randomly replaced by one of the allowed operators. Once again, the new tree is allowed to survive to the next generation only if it is able to capture at least one definition.

Once the GP converges, we expect to have new expressions that would capture some aspects of a definition. The application of this program will allow us to extend our current set of grammar rules by deriving new rules from the above operations. Although we do not expect the GP to learn rules, it will help towards the discovery of new rules which might have been overlooked, and thus helping towards a more complete grammar for definition extraction.

The GP will also allow the flexibility of running this experiment separately for each of the categories of the definitions as identified in section 2.2. This means that the new features being learned will be restricted to one category at a time.

2.3 Combining the Two Experiments

The role of this work is to develop techniques to extract definitions. The two experiments are independent of each other. The GA takes a set of features and assigns a weight to each feature, whereas the GP learns new features through the evolution of the population of extended regular expressions. We can combine the two experiments by migrating the new features learned by the GP to augment the feature set which is used in the GA.

In the final definition extractor one can start by checking whether a given sentence can be confidently classified as a definition by using the features learnt by the GP, possibly giving a preference to over- or under-approximations. One would then run the weighted sum and threshold as learned by the GA based on the features we manually identified and others that the GP may have learned. Clearly the training of the GA would have to be done on a subset of the training set, removing the confidently classified non/definitions. We believe that this approach will improve the quality of the definition identifier.

3 Literature Review

We split this review into two main parts, starting with an overview of published results using rule-based definition extraction, followed by work in grammatical inference which applies GAs and GPs. A final discussion comparing the proposed work to the work reviewed, then follows.

3.1 Rule-Based Definition Extraction

Work carried out on automatic creation of glossaries usually tends to be rule-based, taking into consideration mainly part-of-speech as the main linguistic feature. Park et al. [PBB02] propose a system whereby glossary candidates are presented to an expert in the relevant domain to be approved and made available through an API. In their work, they concentrate on detecting the terms and their glosses rather than full definitions. Glosses present a summary of the meaning, usually giving the full form of an abbreviation or a variant of the term. They also deal only with technical texts, where glosses are normally well structured.

They propose a pipeline architecture using several tools, including POS tagging and morphological analysis, with each tool providing additional annotations. The glossary extraction algorithm first looks at the possible linguistic structure of glossary items found technical texts. They identify POS structures for noun phrases or verbs, which are used to identify possible glosses. These are described as a cascade of finite-state transducers, which are easy to extend and re-use even across different languages. An observation they make is that is difficult to identify glosses by simply looking at POS sequences, since many other non-gloss items would have the same sequence. To overcome this problem, rules are applied to discard certain forms from the candidate set. These include person and place names, special tokens such as URLs, words containing symbols (except for hyphens and dashes) and candidates having more than six words.

Variants are identified, grouped and one is set at the canonical form, others listed as variants (including misspelt items and abbreviations). Finally all glossary candidates are ranked and presented to the expert. In the evaluation of their work, three human experts accepted 228 (76%), 217 (72%), and 203 (68%) out of the top 300 as valid glossary items. The evaluation does not consider missed definition which ranked lower. Inter-annotater agreement is not discussed in this evaluation.

Klavans and Muresan [KM00] propose a system, Finder, to extract definitions from technical, medical texts. The corpora used comprise of consumer-oriented texts, where a medical term is explained in general language words in order to provide a paraphrase. The aim for their system is to be able to extract definitions that can then be fed into a dictionary. Their approach uses NLP techniques to identify definitions and synonyms (which are also considered as definitions in this context). They point out that the structure of definitions might not follow

the *genus et differentia*⁴ model and that the different styles of writing can be a challenge for the automatic identification of definitions.

Definder first identifies candidate phrases by looking for cue-phrases such as “is called a”, “is the term for”, “is defined as”, or a set of punctuation marks which are deemed important for this task (namely :, (,), -). A finite state grammar is then applied to extract the definitions. The system uses part-of-speech and noun phrase chunking to help with the identification process. In order to improve results, the Definder uses statistical information from a grammar analysis module. The authors claim that doing so takes into account the styles for writing of definitions (apposition, relative clauses, anaphora). In this work we see that the automatic identification of definitions is mainly based on the primary identification of certain phrases, and then further filtered through certain rules that reinforce a sentence being a definition (such as its POS structure).

Klavans et al. [KPP03] look at the Internet as a corpus, focusing mainly on large government websites, trying to identify definitions by genus to extract conceptual relations for ontology building. In this task, several problems are identified, including the format of the definitions and the content in which they are present. Definitions on the Internet can be ambiguous, uncertain or incomplete. They are also being derived from heterogeneous document sources. Another problem encountered is that the Internet is a dynamic corpus, and websites could change their information over time. An interesting discussion is presented in how to evaluate a definition extractor, proposing a Gold Standard for such a type of evaluation, based on qualitative evaluation apart from the standard quantitative metrics.

Liu et al. [LCN03] are interested in definition extraction of concepts for learning purposes. Their strategy to assist learning is to present learners with definitions of concepts, and the sub-topics or salient concepts of the original topic. Their system queries search engines with a concept, and the top 100 ranked results are retrieved. In order to discover salient and sub-topics, they look at layout information presented in the html tags for features such as headings, bold and italic. A rule-based approach is then applied to filter out items which are generally also highlighted in webpages (such as company names, URLs, lengthy descriptions). Further filtering is applied through stopword removal, taking frequency into consideration and ranking the proposed salient or sub-topics.

Definition extraction is then attempted for the concepts and sub-concepts identified in the first phase of their work. The identification is carried out through rule-based patterns, (e.g. {concept} {refers to | satisfies}...). Webpages containing definitions are attached to the concept, and presented to the user in ranked order. The ranking is based on how many concepts/definitions are present in a webpage (the more being available, the higher the ranking as it is considered more informative).

Liu et al. also propose a way of dealing with ambiguity, where the term being learned is too generic and may appear in different context (e.g. classification may

⁴ A genus et differentia definition first describes the term explain the broader concept, the *genus*, and then distinguishes from other items in the category by *differentia*

be used in library classification, product classification, data mining classification, etc.). Again, to differentiate between different contexts, the term is allocated a parent topic, and through the use of document layout structure a hierarchical structure of topics is built accordingly. Mutual reinforcement is also used to provide further evidence of the hierarchy built, by further searching for sub-topics under a particular concept. This generally would result in finding information about other salient topics under that same concept, which thus continues to re-ensure that the sub-topics are related and belong to the parent concept.

It is interesting to note that the definition extraction phase is preceded by a phase of concept identification, simplifying the task of definition extraction. However, the search for definitions is carried out on particular concepts, and not all definitions contained in a given text. This might result in definitions of other concepts present in parsed documents being lost. The work also does not make any use of linguistic information, since at their level of processing simple pattern matching on particular keywords is sufficient. The web is a very large corpus and thus can provide many documents containing definitions. However, this can also be a disadvantage affecting the quality of definitions found, similar to the problems encountered by [KPP03] described above (i.e. ambiguous, uncertain or incomplete definitions). This is partially surpassed by providing several resulting definitions, and documents containing more definitions (which might be more authoritative) are presented first. However, quality of definitions is important in any learning phase, and can determine the learner's understanding of a concept.

Storrer and Wellinghoff [SW06] report work on definition classification based on 19 primary verbs, specified in valency frames. These frames indicate what arguments a verb takes, such as object, subject, position and prepositions. This frame can be used to match the structure of a sentence containing one of the specified 19 verbs. Thus definitions are extracted by using the valency frames specified for the defining verbs. The approach presented in the paper is a rule-based expert driven, with all information being provided by human experts (valency frame, definition categorisation). This is possible because they are looking at technical texts, where definitions are well-structured, frequently matching more crisp rules.

Fahmi and Bouma [Fah06] tackle the problem of definition extraction using an incremental approach, starting with individual words, then adding syntactic features etc. They look at the potential definition sentences that fall into our first category (containing the verb to be) from a Dutch corpus of medical articles extracted from Wikipedia. These sentences are manually annotated as definitions, non-definitions and undecided, and this corpus of sentences is used as their training and evaluation data for the experiments carried out. They identify several attributes that could be of importance to the experiments, namely text properties, sentence position, syntactic properties and named entity classes. Learning-based methods are then used to identify which of these features, or combination of, would provide the best results. These feature combinations can also be considered for the experiments described above.

3.2 Grammatical Inference Using GAs and GPs

Identifying grammars for definition extraction is closely related to grammatical inference — the use of machine learning techniques to learn grammars from data. GAs, and less frequently GPs, are two such techniques which have been applied to grammatical inference.

Genetic Algorithms

Lankhorst [Lan94] describes a genetic algorithm used to infer context-free grammars from positive and negative examples of a language. The schema theorem states that schemas occurring in the higher scoring individuals will tend to occur more frequently in following generations. This feature, also referred to as the *building blocks hypothesis*, is the motivation for applying GAs to grammatical inference. New, possibly better performing, grammar rules may be discovered by combining parts of different grammar rules.

A discussion is presented on the choice of gene representation, between a binary representation and a high-level representation. It is argued that a bit string can represent many more schemata than higher level representations, yielding to a better coverage of the search space. A bit representation is chosen, with the lower order bits encoding grammar rules on the right hand side of the rule, whereas higher order bits are encoding the left-hand side symbol.

Selection is based on a stochastic universal sampling algorithm, that helps to prevent premature convergence by ‘holding back’ super individuals from taking over the population within a few generations. The best individual is also always allowed to survive to the next generation. Mutation allows for a bit in a chromosome to be mutated. However, this operation is given a low probability so as not to change the population randomly. Reproduction is effected by the schema chosen. The crossover point is influenced by how the representation of the rule is expressed. Lankhorst chooses a two-point crossover system, thus allowing right-hand side rules to crossover more easily.

The grammars that Lankhorst is learning with the GA are aimed to classify positive and negative examples over a language correctly. Thus the fitness function that is considered takes into consideration the correct classification of positive and negative examples. For various sets of languages, such as matching bracketing and ‘0*1*0*1*’, the fitness function allows the population to converge into a solution within reasonable results. However, for a micro-NL language, the fitness function did not result in a correct grammar. A further adjustment to the fitness function also allows correct classifications of substrings to influence the fitness score. This modified fitness function allows the GA for a micro-NL grammar to converge correctly.

This work provides an interesting insight to different techniques of how the fitness of an individual should be calculated. Different fitness functions will output different results and it is important to explore alternatives. In our case we would like to take into consideration which measures should be given more importance to (e.g. recall, precision or F-measure). The representation of an individual is

also important, Lankhorst having selected a bit-representation over a higher-level representation such as trees. Tree representations can be easily applied with Genetic Programming, and yet, exactly the same principles/discussions apply with respect to the fitness function.

Losee [Los96] applies a GA to learn the syntactic rules and tags with a direct application in document retrieval. The system parses just over 100 abstracts, all about the same general topic, subdivided into 5 different subtopics. The task of the system is to retrieve the documents in the ideal ranking order according to the search term used. The GA is applied to learn syntactic rules and tags, and thus provides linguistic meaning to both documents and terms. As a fitness function, the GA uses a weighted function of the resulting ranking of the document and the average maximum parse length.

Belz and Eskikaya [BE98] attempt grammatical induction from positive data sets in the field of phonotactics. The grammar is represented with finite state automata, and looks at the use of GA for this task. In this paper two results are produced, one for German syllables and the other for Russian bisyllabic words. The GA is described in detail, including the type of methods selected, and chromosome representation. The GA used is a fine-grained one, where individuals are on a 2-D grid of fixed size and are only allowed to mate with one of their neighbours (this implementation is referred to as a torus).

They argue that an important issue is the representation used for individuals. They present two alternatives: (1) production rules of the form $s_1 \rightarrow as_2$ and $s_1 \rightarrow a$ (where s is a non-terminal symbol and a is a terminal symbol), or (2) a state transition matrix. They argue that production rules produce more fine grained genotype representations since the terminals and non-terminals can be represented individually. On the other hand, state transition matrices can be only seen as a whole, each represented by a single cell. The final representation chosen for this work is that of transition matrices. Each individual represents a possible transition matrix which in turn represents the grammar being induced. In order for genetic algorithms to be used, the matrix is ‘flattened’ into one string (one row after the other).

The chosen representation has direct implications on the rest of the GA operations. Crossover and mutation cannot be carried out in the traditional sense of GAs, and certain knowledge must be present in the GA so as to maintain a sound structure of this flat matrix. Belz and Eskikaya seem not to have considered GPs and tree representation for their problem. Such consideration could have provided an interesting alternative in their work.

Keller and Lutz [KL97] attempt learning Context Free Grammars through the use of GAs, by learning probabilities to all possible grammar rules. The initial population of the GA is made of all possible combinations of terminals and non-terminals of the form $A \rightarrow BC$ and $A \rightarrow a$, where A, B and C are non-terminals and a is a terminal. This guarantees that, although the grammar is a large one, it is finite. There is also no loss of generality, as all possible rules are present in the initial grammar (including those that will not be part of the final solution).

The type of GA chosen for this work also uses a 2-D grid representation in torus formation, where mating occurs only with immediate neighbours. Each individual is encoded as a set of weights, each weight relevant to one parameter. The weight is represented as an n -bit block, and an individual can be viewed as consisting of M blocks of n -bits, where M is the number of all possible rules. Since the grammar is considerably larger than the final solution, Keller and Lutz try to give more importance to 0-probability assignment. Thus, the initial bit/s of the n -bit block is seen as a “binary-switch” as to whether the rest of the bits should be taken into consideration or not.

As for crossover, Keller and Lutz achieved better performance by using a novel genetic operator which they call *and-or crossover* than by the classical crossover operation. The and-or crossover looks at the parents bit by bit, with one child taking the bits produced by the and operator (conservative), and the other child takes the bits produced by the or operator (liberal). In our proposed work, this idea of the bit-by-bit crossover can be adapted slightly to use functions such as minimum, average and maximum to develop new children.

Spasić et al. [SNA04] aim at classifying biomedical terms into specific classes, which represent concepts from an ontology in the biomedicine domain. In order to derive the possible class of a term, they look at the surrounding context of that term. This context is learned through data mining, extracting the contextual patterns surrounding the terms. The patterns contain morph-syntactic and terminological information and are represented as generalised regular expressions. Each contextual pattern is then given a value indicating its statistical relevance. They use this to remove the top and bottom ranked features since they are considered too general or too rare to play a role in term classification. Class selection of a term is then learned using a Genetic Algorithm. For a particular class, the GA tries to learn which of the contextual patterns are relevant. Each individual in the GA is a subset of contextual patterns and its fitness corresponds to the precision and recall of using these patterns on the training data. Eventually the GA learns a good subset of features which can be used to identify terms in that class.

Genetic Programming

Smith and Witten [SW95] propose a GP that adapts a population of hypothesis grammars towards a more effective model of language structure. They discuss grammatical inference using statistical methods, and the problems encountered in their work. They point out that probabilistic n -gram models allow frequent, well-formed expressions to statistically overwhelm infrequent ungrammatical expressions. There is also the problem with allowing probability for unseen data. The ‘zero-frequency problem’ entails assigning a small probability to all unseen data, resulting in both ungrammatical n -grams becoming as probable as unseen grammatical ones.

The population is represented as LISP AND-OR S-expressions. Initial experiments showed that certain constraints were required in order for the GP to evolve. These constraints included a maximum depth for nesting and a grammar-

generator to allow the GP to evolve towards more suitable grammars. With these constraints in place, the GP evolves simple grammars even within 2 generations, forming simple sentences such as ‘the dog saw a cat’. However, the GP is left to run over more generations to achieve a broader exploration of the search space and hopefully result in a more efficient grammar. Although the results seem positive, there is no comparison to other statistical techniques mentioned that attempt grammatical inference.

3.3 Discussion

Work carried out in definition extraction shows that although it is possible to achieve a good basis for a grammar through manual observation, this task requires specialized linguistic understanding of grammatical features present in definitions. Ideally, a filtering or ranking mechanism is used to over and above these techniques to further improve the results. Work by Fahmi and Bouma [Fah06] moves towards this direction. In our proposed experiments we introduce the concept of ranking through the use of the GA, learning weights assigned to grammars indicating the certainty of whether the captured sentences are actually definitions.

The GP, and in particular the gene representation, is a crucial point which comes out clearly in the work reviewed. There is a lack of knowledge in the area of GPs, and this is visible in the various attempts of describing a grammar as a linear structure to work with GAs rather than taking advantage of a GP’s capability to handle tree representation. Smith and Witten [SW95] overcome representation issues by using grammars as LISP S-expressions. In our work, the resulting grammar that will be produced will be used within *lxtransduce* [Tob05]. Thus any type of representation chosen will be translated to the XML format accepted by *lxtransduce*.

It is also clear that the fitness function will determine the success of the experiments. Since we have to our availability manually annotated definitions, precision and recall could be used as part of the fitness function. However, different tests could be carried out to determine what should comprise the fitness evaluation of the population.

4 Conclusion

Attempts at definition extraction have focused mainly on rule-based approaches, with some later work improving results by introducing statistical analysis as a filtering step. Our proposal introduces an element of grammar learning for the set of definitional sentences, influenced by the work carried out in grammatical inference. We will also use the weights produced by the GA as part of the filtering and ranking process, as evidence to what should be classified as a definition. The proposal takes a novel approach in combining GAs and GPs for natural language processing. A quantitative evaluation of these techniques will compare

the results achieved to the work carried out so far in definition extraction. The results should be of interest not only to the natural language task of extracting definitions, but also to the machine learning task of combining GAs with GPs.

References

- [BE98] Anja Belz and Berkan Eskikaya. A genetic algorithm for finite state automata induction with an application to phonotactics. In *Proceedings of the ESSLLI-98 Workshop on Automated Acquisition of Syntax and Parsing*, 1998.
- [BR07] Claudia Borg and Mike Rosner. Language Technologies for an eLearning Scenario. In *Proceedings of the Computer Science Annual Workshop, CSAW07, Malta*, 2007.
- [Fah06] Learning to Identify Definitions using Syntactic Features. In *Workshop of Learning Structured Information in Natural Language Applications, EACL, Italy*, 2006.
- [Gol89] David E. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley, Reading, MA, 1989.
- [Hol75] John H. Holland. *Adaptation in Natural and Artificial Systems*. University of Michigan Press, Ann Arbor, 1975.
- [KL97] Bill Keller and Rudi Lutz. Learning Stochastic Context-Free Grammars from Corpora Using a Genetic Algorithm. In *Workshop on Automata Induction Grammatical Inference and Language Acquisition (ICML-97)*, Nashville, Tennessee, 1997.
- [KM00] Judith L. Klavans and Smaranda Muresan. DEFINDER: Rule-Based Methods for the Extraction of Medical Terminology and their Associated Definitions from On-line Text. In *Proceedings of the American Medical Informatics Association (AMIA) Symposium*, 2000.
- [Koz92] John R. Koza. *Genetic Programming: On the Programming of Computers by means of Natural Selection*. MIT Press, Cambridge, MA, 1992.
- [KPP03] Judith L. Klavans, Samuel Popper, and Rebecca Passonneau. Tackling the internet glossary glut: Automatic extraction and evaluation of genus phrases. In *SIGIR'03 Workshop on Semantic Web*, 2003.
- [Lan94] Marc M. Lankhorst. Breeding grammars: Grammatical inference with a genetic algorithm. Technical Report CS-R9401, Department of Computer Science, University of Gronigen, PO Box 800, 9700 AV Groningen, The Netherlands, 1994.
- [LCN03] Bing Liu, Chee W. Chin, and Hwee T. Ng. Mining Topic-Specific Concepts and Definitions on the Web. In *Proceedings of the Twelfth International World Wide Web Conference (WWW'03)*, 2003.
- [Los96] Robert M. Losee. Learning Syntactic Rules and Tags with Genetic Algorithms for Information Retrieval and Filtering: An Empirical Basis for Grammatical Rules. In *Information Processing and Management*, volume 32, 1996.
- [MLS07] Paola Monachesi, Lothar Lemnitzer, and Kiril Simov. Language Technology for eLearning. In *First European Conference on Technology Enhanced Learning*, 2007.
- [PBB02] Youngja Park, Roy J Byrd, and Branimir K Boguraev. Automatic glossary extraction: beyond terminology identification. In *Proceedings of the 19th international conference on Computational linguistics*, pages 1–7, Morristown, NJ, USA, 2002. Association for Computational Linguistics.

- [SNA04] Irena Spasić, Goran Nenadić, and Sophia Ananiadou. Learning to Classify Biomedical Terms through Literature Mining and Genetic Algorithms. In *Intelligent Data Engineering and Automated Learning (IDEAL 2004)*, volume 3177 of *LNCS*, pages 345–351. Springer-Verlag, 2004.
- [SW95] Tony C. Smith and Ian H. Witten. A Genetic Algorithm for the Induction of Natural Language Grammars. In *Proceedings IJCAI-95 Workshop on New Approaches to Learning for Natural Language Processing, Canada*, pages 17–24, 1995.
- [SW06] Angelika Storrer and Sandra Wellinghoff. Automated detection and annotation of term definitions in german text corpora. In *LREC*, 2006.
- [Tob05] Richard Tobin. Lxtransduce a replacement for fsmatch. Technical report, University of Edinburgh, 2005.