

Grammatical Inference in Software Engineering: an Overview of the State of the Art

Andrew Stevenson and James R. Cordy

Queen's University, Kingston ON, Canada
{andrews, cordy}@cs.queensu.ca

Abstract. Grammatical inference – used successfully in a variety of fields such as pattern recognition, computational biology and natural language processing – is the process of automatically inferring a grammar by examining the sentences of an unknown language. Software engineering can also benefit from grammatical inference. Unlike the aforementioned fields, which use grammars as a convenient tool to model naturally occurring patterns, software engineering treats grammars as first-class objects typically created and maintained for a specific purpose by human designers. We introduce the theory of grammatical inference and review the state of the art as it relates to software engineering.

Keywords: grammatical inference, software engineering, grammar induction

1 Introduction

The human brain is extremely adept at seeing patterns by generalizing from specific examples, a process known as inductive reasoning. This is precisely the idea behind grammatical induction, also known as grammatical inference, where the specific examples are sentences and the patterns are grammars. Grammatical inference is the process of identifying an unknown language by examining examples of sentences in that language. Specifically, the input to the process is a set of strings and the output is a grammar.

The main challenge of identifying a language of infinite cardinality from a finite set of examples is knowing when to generalize and when to specialize. Most inference techniques begin with the given sample strings and make a series of generalizations from them. These generalizations are typically accomplished by some form of state-merging (in finite automata), or non-terminal merging (in context-free grammars).

Grammatical inference techniques are used to solve practical problems in a variety of different fields: pattern recognition, computational biology, natural language processing and acquisition, programming language design, data mining, and machine learning. Software engineering, in particular software language engineering, is uniquely qualified to benefit because it treats grammars as first-class objects with an intrinsic value rather than simply as a convenient mechanism to model patterns in some other subject of interest.

Historically there have been two main groups of contributors to the field of grammatical inference: theorists and empiricists. Theorists consider language classes and learning models of varying expressiveness and power, attempting to firm up the boundaries of what is learnable and how efficiently it can be learned, whereas empiricists start with a practical problem and, by solving it, find that they have made a contribution to grammatical inference research.

Grammatical inference is, intuitively as well as provably, a difficult problem to solve. The precise difficulty of a particular inference problem is dictated by two things: the complexity of the target language and the information available to the inference algorithm about the target language. Naturally, simpler languages and more information both lead to easier inference problems. Most of the theoretical literature in this field investigates some specific combination of language class and learning model, and presents results for that combination.

In Section 2 we describe different learning models along with the type of information they make available to the inference algorithm. In Section 3 we explore the learnability, decidability, and computational complexity of different learning models applied to language classes of interest in software engineering: finite state machines and context-free grammars. Section 4 discusses the relationship between theoretical and empirical approaches, and gives several practical examples of grammatical inference in software engineering. In Section 5 we list the related surveys, bibliographies, and commentaries on the field of grammatical inference and briefly mention the emphasis of each. Finally, in Section 6 we discuss the main challenges currently facing software engineers trying to adopt grammatical inference techniques, and suggest future research directions to address these challenges.

2 Learning Models

The type of learning model used by an inference method is fundamental when investigating the theoretical limitations of an inference problem. This section covers the main learning models used in grammatical inference and discusses their strengths and weaknesses.

Section 2.1 describes *identification in the limit*, a learning model which allows the inference algorithm to converge on the target grammar given a sufficiently large quantity of sample strings. Section 2.2 introduces a teacher who knows the target language and can answer particular types of queries from the learner. This learning model is, in many cases, more powerful than learning from sample strings alone. Finally, Section 2.3 discusses the PAC learning model, an elegant method that attempts to find an optimal compromise between accuracy and certainty. Different aspects of these learning models can be combined and should not be thought of as mutually exclusive.

2.1 Identification in the limit

The field of grammatical inference began in earnest with E.M. Gold’s 1967 paper, titled “Language Identification in the Limit” [24]. This learning model provides

the inference algorithm with a sequence of strings one at a time, collectively known as a presentation. There are two types of presentation: positive presentation, where the strings in the sequence are in the target language; and complete presentation, where the sequence also contains strings that are not in the target language and are marked as such. After seeing each string the inference algorithm can hypothesize a new grammar that satisfies all of the strings seen so far, i.e. a grammar that generates all the positive examples and none of the negative examples. The term “information” is often used synonymously with “presentation” (e.g. positive information and positive presentation mean the same thing).

The more samples that are presented to the inference algorithm the better it can approximate the target language, until eventually it will converge on the target language exactly. Gold showed that an inference algorithm can identify an unknown language in the limit from complete information in a finite number of steps. However, the inference algorithm will not know when it has correctly identified the language because there is always the possibility the next sample it sees will invalidate its latest hypothesis.

Positive information alone is much less powerful, and Gold showed that any superfinite class of languages cannot be identified in the limit from positive presentation. A superfinite class of languages is a class that contains all finite languages and at least one infinite language. The regular languages are a superfinite class, indicating that even the simplest language class in Chomsky’s hierarchy of languages is not learnable from positive information alone.

There has been much research devoted to learning from positive information because the availability of negative examples is rare in practice. However, the difficulty of learning from positive data is in the risk of overgeneralization, learning a language strictly larger than the target language. Angluin offers a means to avoid overgeneralization via “tell-tales”, a unique set of strings that distinguish a language from other languages in its family [2]. She states conditions for the language family that, if true, guarantee that if the tell-tale strings are included in the positive presentation seen so far by the inference algorithm then it can be sure its current guess is not an overgeneralization.

2.2 Teacher and Queries

This learning model is similar in spirit to the game “twenty questions” and uses a teacher, also called an oracle, who knows the target language and answers queries from the inference algorithm. In practice, the teacher is often a human who knows the target language and aids the inference algorithm, but in theory can be any process hidden from the inference algorithm that can answer particular types of questions. Angluin describes six types of queries that can be asked of the teacher, two of which have a significant impact on language learning: membership and equivalence [6]. A teacher that answers both membership and equivalence queries is said to be a *minimally adequate teacher* because she is sufficient to help identify DFAs in polynomial time without requiring any examples from the target language [5].

For a membership query, the inference algorithm presents a string to the teacher who responds with “yes” if the string is in the language or “no” if it is not. Likewise for an equivalence query, the inference algorithm presents a grammar hypothesis to the teacher who answers “yes” or “no” if the guess is equivalent to the target grammar or not. In the case when the teacher answers “no” she also provides a counter-example, a string from the symmetric difference of the target language and the guessed language, allowing the inference algorithm to zero in on the target grammar. The symmetric difference of two sets A and B are the elements in either A or B but not both: $A \oplus B = (A \cup B) - (A \cap B)$.

Queries provide an alternate means to measure the learnability of a class of languages. They can be used on their own or in conjunction with a presentation of samples, either positive or complete, to augment the abilities of the learner. Section 3 discusses how learning with queries differs in difficulty from learning in the limit for various language classes.

2.3 PAC Learning

In 1984 Valiant proposed the Probably Approximately Correct (PAC) learning model [55]. This model has elements of both identification in the limit and learning from an oracle, but differs because it doesn’t guarantee exact identification with certainty. As its name implies, PAC learning measures the correctness of its result by two user-defined parameters, ϵ and δ , representing accuracy and confidence respectively. This learning model is quite general and thus uses different terminology than typically found in formal languages, but of course applies just as well to grammatical inference. The goal is still to learn a “concept” (grammar) from a set of “examples of a concept” (strings).

Valiant assumes there exists a (possibly unknown) distribution D over the examples of a target concept that represent how likely they are to naturally occur, and makes available to the inference algorithm a procedure that returns these examples according to this distribution. As with Gold’s identification in the limit, PAC learning incrementally approaches the target concept with more accurate guesses over time.

A metric is proposed to measure the distance between two concepts, defined as the sum of probabilities $D(w)$ for all w in the symmetric difference of $L(G)$ and $L(G')$. In Figure 1, the lightly shaded regions represent the symmetric difference between $L(G)$ and $L(G')$. The area of this region decreases as the distance between the two concepts decreases. In the case of grammatical inference, these two concepts refer to the target grammar and the inference algorithm’s current guess.

The PAC learning model’s criteria for a successful inference algorithm is one that can confidently (i.e. with probability at least $1 - \delta$) guess a concept with high accuracy (i.e. distance to the target concept is less than ϵ). Valiant demonstrates the PAC learning model with a polynomial time algorithm that approximates bounded conjunctive normal form (k-CNF) and monotone disjunctive normal form (DNF) expressions using just the positive presentation from D and a membership oracle.

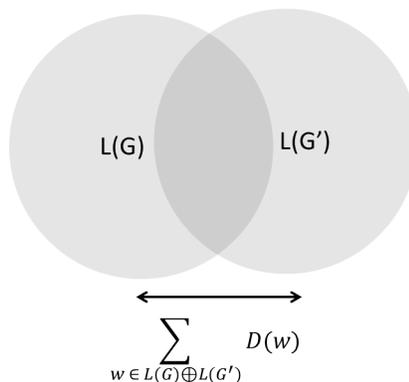


Fig. 1. The PAC-learning measure of distance between two language concepts

The novelty and uniqueness of Valiant’s model intrigued the grammatical inference community, but negative and NP-hardness equivalence results (e.g. [33, 49]) dampened enthusiasm for PAC learning. Many feel Valiant’s stipulation that the algorithm must learn polynomially under *all* distributions is too stringent to be practical since the learnability of many apparently simple concept classes are either known to be NP-hard, or at least not known to be polynomially learnable for all distributions.

Li and Vitanyi propose a modification to the PAC learning model that only considers simple distributions [43]. These distributions return simple examples with high probability and complex examples with low probability, where simplicity is measured by Kolmogorov complexity. Intuition is that simple examples speed learning. This is corroborated by instances of concepts given by the authors that are polynomially learnable under simple distributions but not known to be polynomially learnable under Valiant’s more general distribution assumptions.

Despite the learnability improvements that simple PAC learning offers, the PAC learning model has attracted little interest from grammatical inference researchers in recent years. Identification in the limit and query-based learning models remain far more prevalent, with newer models such as neural networks and genetic algorithms also garnering interest.

3 Complexity

A significant portion of the grammatical inference literature is dedicated to an analysis of its complexity and difficulty, with results typically stated for a specific grammar class or type. The broadest form of result is simply whether a language class can be learned or not, while other results consider learning in polynomial time, learning the simplest grammar for the target language, or identifying the target language with a particular probability. Table 1 outlines the complexity results for different language classes and learning models.

Language Class	Presentation		Queries		
	Complete	Positive	Membership Only	Equivalence Only	Both
Finite	Identifiable in the limit [24]	Identifiable in the limit [24]			
k-reversible automata		Polynomial [4]			
Strictly deterministic automata		Identifiable in the limit [60]			
Superfinite	Identifiable in the limit [24]	Not identifiable in the limit [24]			
Regular	Finding the minimum state DFA is NP-hard [25]		Polynomial for representative sample [3]	No polynomial algorithm [7]	Polynomial [5]
			Polynomial [15]		
Reversible context-free		Identifiable in the limit with structured strings [52]			
Noncounting context-free		Identifiable with structured strings [16]			
Very simple		Polynomial identifiable in the limit [59]		Polynomial [59]	
Structurally reversible context-free					Polynomial [11]
Simple deterministic					Polynomial [28]
Context-free			As hard as inverting RSA [8]		Polynomial with structured strings [51]

Table 1. Learnability and complexity results for various language classes using different learning models

Gold showed that a large class of languages can be identified in the limit from complete information including the regular, context-free, context-sensitive, and primitive recursive classes. This identification can be accomplished by a brute-force style of technique called *identification by enumeration* where, as each new example is presented, the possible grammars are enumerated until one is found that satisfies the presentation seen so far. By contrast, positive information alone cannot identify the aforementioned classes in the limit, nor any other superfinite class [24]. The subsequent sections describe the two easiest grammar classes to infer from the Chomsky hierarchy: regular grammars and context-free grammars. Very little research has been attempted on the inference of more powerful grammar classes such as context-sensitive and unrestricted grammars, so they are omitted from this overview.

3.1 Deterministic Finite Automata

For any non-trivial language, multiple different grammars can be constructed to generate it. Likewise, there can exist DFAs that differ in their size but are equivalent in the sense that they accept the same language. When inferring a DFA from examples, it is naturally desirable to find the smallest DFA that accepts the target language. There exists only one such minimal DFA for a given language, known as the *canonical* DFA acceptor for that language. Despite the strong identification power of complete information, finding the minimal DFA that accepts an unknown regular language from a finite set of positive and negative samples is NP-complete [25].

Early claims of polynomial-time inference algorithms use the number of states in the target language's canonical acceptor as the input size. With this criteria, Angluin gives negative results for the polynomial identification of regular languages using membership queries only [3] or equivalence queries only [7]. However, if the membership oracle is augmented with a *representative sample* of positive data, a set of strings that exercise all the live transitions in the target language's canonical acceptor, then it is possible to identify a regular language in polynomial time [3]. By combining the power of both membership and equivalence queries, a regular language can be identified in polynomial time even without a single positive example in the unknown language [5]. Her proposed algorithm runs in time polynomial to the number of states in the minimum DFA and the longest counter-example provided by the equivalence oracle.

Several algorithms have been developed for the inference of DFAs from examples. These algorithms generally start by building an augmented prefix tree acceptor from positive and negative samples, then perform a series of state merges until all valid merges are exhausted. Each state merge has the effect of generalizing the language accepted by the DFA. The algorithms differ by how they select the next states to merge, constraints on the input samples, and whether or not they guarantee the inference of a minimal DFA.

An early state-merging algorithm is described by Trakhtenbrot and Barzdin that infers a minimal DFA in polynomial time, but requires that all strings up to a certain length are labeled as positive or negative [54]. The *regular positive*

and negative inference (RPNI) algorithm also finds the canonical acceptor but allows an incomplete labeling of positive and negative examples which is more common in practice [48]. RPNI does, however, require the positive examples contain a characteristic set with respect to the target acceptor. A *characteristic set* is a finite set of positive examples $S \subset L(A)$ such that there is no other smaller automata A' where $S \subset L(A') \subset L(A)$. Lang provides convincing empirical evidence that his *exbar* algorithm out-performs comparable algorithms, and represents the state of the art in minimal DFA inference [38].

The algorithms discussed so far are guaranteed to infer a minimal DFA for the given examples, but *evidence driven state merging* (EDSM) algorithms relax this requirement for better scalability and performance. The order that states in a prefix tree are merged has a significant impact on an algorithm's performance because each merge restricts possible future merges. Bad merge decisions cause a lot of backtracking that can be avoided with smarter merge decisions. EDSM algorithms are so named because they use evidence from the merge candidates to determine a merge that is likely to be a good generalization, such as the heuristic proposed by Rodney Price in the first EDSM algorithm [39] and a winner of the Abbadingo Learning Competition. Differences in EDSM algorithms come down to the search heuristic used to select merges, and several have been tried such as beam search [38], stochastic search and the *self-adaptive greedy estimate* (SAGE) algorithm [32]. These search heuristics are comparable in performance and are the best known inference algorithms for large or complex DFAs.

3.2 Context-free grammars

Polynomial-time algorithms to learn higher grammar classes have also been investigated, in particular for context-free grammars. Identifying context-free grammars in polynomial time is considerably more difficult than for DFAs, so most polynomial results in the literature either learn a strict subset of context-free grammars, use structured strings as input, or both. Unlike DFA inference, there is currently no known polynomial algorithm to identify a general context-free language from positive and negative samples.

Angluin and Kharitonov give a hardness result that applies to all context-free languages: constructing a polynomial-time learning algorithm for context-free grammars using membership queries only is computationally equivalent to cracking well-known cryptographic systems, such as RSA inversion [8].

Anecdotally, it appears a fruitful method to find polynomial-time learning algorithms for context-free languages from positive samples is to adapt corresponding algorithms from DFA inference, with the added stipulation that the sample strings be structured. A structured string is a string along with its unlabelled derivation tree, or equivalently a string with nested brackets to denote the shape of its derivation tree. Sakakibara has shown this method effective by adapting Angluin's results for learning DFAs by a minimally adequate teacher [5] and learning reversible automata [4] to context-free variants with structured strings [51, 52].

Clark et al. have devised a polynomial algorithm for the inference of languages that exhibit two special characteristics: the finite context property and the finite kernel property [15]. These properties are exhibited by all regular languages, many context-free languages, and some context-sensitive languages. The algorithm is based on positive data and a membership oracle. More recently, Clark has extended Angluin’s result [5] of learning regular languages with membership and equivalence queries to a larger subclass of context-free languages [14].

Despite the absence of a general efficient context-free inference algorithm, many researchers have developed heuristics that provide relatively good performance and accuracy by sacrificing exact identification in all cases. We describe several such approaches related to software engineering in Section 4.

4 Applications in Software Engineering

Grammatical inference has its roots in a variety of separate fields, a testament to its wide applicability. Implementors of grammatical inference applications often have an unfair advantage over purely theoretical GI research because theorists must restrict themselves to inferring abstract machines (DFAs, context-free grammars, transducers, etc.) making no additional assumptions about the underlying structure of the data. Empiricists, on the other hand, can make many more assumptions about the structure of their data because their inference problem is limited to their particular domain.

Researchers attempting to solve a practical inference problem will usually develop their own custom solution, taking advantage of structural assumptions about their data. Often this additional domain knowledge is sufficient to overcome inference problems that theorists have proved impossible or infeasible with the same techniques in a general environment. The applications described in the following sections use grammatical inference techniques, but rarely result from applying a purely theoretical result to a practical problem.

4.1 Inference of General Purpose Programming Languages

Programming language design is an obvious area to benefit from grammatical inference because grammars themselves are first-class objects. Programming languages almost universally employ context-free, non-stochastic grammars to parse a program, which narrows the possible inference approaches considerably when looking for an inductive solution. When discussing the inference of programming language grammars here, the terms “sample” and “example” refer to instances of computer programs written in the target programming language.

Crespi-Reghizzi et al. suggest an interactive system to semi-automatically generate a programming language grammar from program samples [17]. This system relies heavily on the language designer to help the algorithm converge on the target language by asking for appropriate positive and negative examples. Every time the learning algorithm conjectures a new grammar, it outputs all sentences for that grammar up to a certain length. If the conjectured grammar

is too large, there will be sentences in the output that don't belong and the designer marks them as such. If the conjectured grammar is too small, there will be sentences missing from the output and the designer is expected to provide them. The designer's corrections are fed back into the algorithm which corrects the grammar and outputs a new conjecture, and the process repeats until the target grammar is obtained.

Another system is proposed by Dubey et al. to infer a context-free grammar from positive samples for a programming language dialect when the standard language grammar is already known [21]. Their algorithm requires the non-terminals in the dialect grammar to be unchanged from the standard grammar, but allows for the terminals and production rules to be extended in the dialect grammar (i.e. new keywords can be added in the dialect along with their associated grammar rules). Their approach has the advantage of being fully automated so the designer simply needs to provide the dialect program samples and the standard language grammar. However, like many current CFG inference techniques, a heuristic is used which cannot guarantee the output grammar converges exactly to the target grammar.

4.2 Inference of Domain Specific Languages

Domain specific languages (DSLs) are languages whose syntax and notation are customized for a specific problem domain, and are often more expressive and declarative compared to general purpose languages. DSLs are intended to be used, and possibly designed, by domain experts who do not necessarily have a strong computer science background. Grammatical inference allows the creation of a grammar for a DSL by only requiring positive (and possibly negative) program samples by the designer.

Črepinšek et al. propose a genetic approach to infer grammars for small DSLs using positive and negative samples [56]. They combine a set of grammar production rules into a *chromosome* representing a complete grammar, then apply crossover and mutation genetic operators that modify a population of chromosomes for the next generation. They use a fitness function that reflects the goal of having the target grammar accept all positive samples and reject all negative samples. Since a single random mutation is more likely to produce a grammar that rejects both positive and negative samples, the authors found that testing a chromosome on only positive samples converges more quickly to the target grammar than testing it on negative samples. Therefore, they chose a fitness value proportional to the total length (in tokens) of the positive samples that can be parsed by a chromosome. Negative samples, used to control overgeneralization, are only included in the fitness value if all positive samples are successfully parsed.

This genetic approach has been shown to accurately infer small DSLs [56], including one discussed by Javed et al. to validate UML class diagrams from use cases [29]. Javed et al. express UML class diagrams in a custom DSL and require a domain expert to provide positive and negative use cases written in that DSL. The system validates these use cases against the given UML diagrams

and reports feedback to the user, who can use that feedback to change the UML diagrams to improve use case coverage. In this situation the computer is providing valuable context and information to the human user who is making the important generalization and specialization decisions for the grammar, but in theory UML diagrams can be synthesized entirely from the use case descriptions given a sufficiently powerful grammar inference engine.

Javed et al. extend their genetic algorithm by learning from positive samples only by using beam search and Minimum Description Length (MDL) heuristics [40] in place of negative examples to control overgeneralization of the conjectured grammar [31]. The idea here is to find the simplest grammar at each step and incrementally approach the target grammar. MDL is used as a measure of grammar simplicity, and beam search is used to more efficiently search the solution space of possible grammars. One disadvantage of this approach is it requires the positive samples to be presented in a particular order, from simplest to most complex, which allows the learning algorithm to encode the incremental differences from the samples into the target grammar. The authors' subsequent effort into a grammar inference tool for DSLs, called *MAGIc*, eliminates this need for an order-specific presentation of samples by updating the grammar based on the difference between successive (arbitrary) samples [46, 27]. This frees the designer from worrying about the particular order to present their DSL samples to the learning algorithm. Hrnčič et al. demonstrate how *MAGIc* can be adapted to infer DSLs embedded in a general purpose language (GPL) given the GPL's grammar [26]. The GPL's grammar rules are included in the chromosome, but frozen so they cannot mutate. Learning, therefore, occurs strictly on the DSL syntax and the locations in the GPL grammar where the embedded DSL is allowed.

The inference of DSLs can make it easier for non-programmer domain experts to write their own domain-specific languages by simply providing examples of their DSL programs. It can also be used in the migration or maintenance of legacy software whose grammar definitions are lost or unavailable.

4.3 Inference of Graph Grammars and Visual Languages

Unlike one-dimensional strings whose elements are connected linearly, visual languages and graphs are connected in two or more dimensions allowing for arbitrary proximity between elements. Graph grammars define a language of valid graphs by a set of production rules with subgraphs instead of strings on the right-hand side.

Fürst et al. propose a graph grammar inference algorithm based on positive and negative graph samples [23]. The algorithm starts with a grammar that produces exactly the set of positive samples then incrementally generalizes towards a smaller grammar representation, a strategy similar to typical DFA inference algorithms which build a prefix tree acceptor then generalize by merging states. The authors demonstrate their inference algorithm with a flowchart example and a hydrocarbon example, making a convincing case for its applicability to

software engineering tasks such as metamodel inference and reverse engineering visual languages.

Another graph grammar inference algorithm is proposed by Ates et al. which repeatedly finds and compresses overlapping identical subgraphs to a single non-terminal node [9]. This system uses only positive samples during the inference process, but validates the resulting grammar by ensuring all the graphs in the training set are parsable and other graphs which are close to but distinct from the training graphs are not parsable. Ates et al. demonstrate their algorithm with two practical inferences: one for the structure of a programming language and one for the structure of an XML data source.

Kong et al. use graph grammar induction to automatically translate a webpage designed for desktop displays into a webpage designed for mobile displays [35]. The inference performed is similar to the aforementioned proposed by Ates et al. [9] because they both use the Spatial Graph Grammar (SGG) formalism and subgraph compression. The induction algorithm consumes webpages, or more accurately their DOM trees, to produce a graph grammar. After a human has verified this grammar it is used to parse a webpage, and the resulting parse is used to segment the webpage into semantically related subpages suitable for display on mobile devices.

Graph grammar inference algorithms are less common than their text-based counterparts, but provide a powerful mechanism to infer patterns in complex structures. Parsing graphs is NP-hard in general, causing these algorithms to be more computationally expensive than inference from text. Most graph grammar learners overcome this complexity by restricting their graph expressiveness or employing search and parse heuristics to achieve a polynomial runtime.

4.4 Other uses in Software Engineering

Section 3 describes the difficulty inferring various language classes from positive samples alone, and in particular that only finite languages can be identified in the limit from positive samples [24]. The SEQUITUR algorithm, developed by Nevill-Manning and Witten, is designed to take a single string (long but finite) and produce a context-free grammar that reflects repetitions and hierarchical structure contained in that string [47]. This differs from typical grammar inference algorithms because it does not generalize. Data compression is an obvious use of this algorithm, but it has found other uses in software engineering. For example, Larus uses the SEQUITUR algorithm to concisely represent a program's entire runtime control flow and uses this dynamic control flow information to identify heavily executed paths in the program to focus performance and compiler optimization efforts [41]. It can also be used on the available positive samples as a first step in a generalizing context-free grammar inference algorithm, such as in [46] to seed an initial population of grammars for a genetic approach.

Ahmed Memon proposes using grammatical inference in log files to identify anomalous activity [45]. He treats the contents of log files in a system running normally as a specific language, and any erroneous or anomalous activities reported in the log file are therefore not part of this language. Memon trains

a grammar from positive log file samples of a system running normally, then parses subsequent log file entries using this grammar to identify anomalous activity. The inference procedure depends on knowledge of the domain, specifically sixteen text patterns that appear in typical log files: dates, times, IP addresses, session IDs, etc. Once these patterns are identified and normalized, a custom non-terminal merging algorithm is used to generalize the log file grammar.

Another recent use for grammatical inference is in the area of model-driven engineering. The relationship between a grammar and the strings it accepts is analogous to the relationship between a metamodel and the instance models it accepts. Javed et al. describe a method to use grammar inference techniques on a set of instance models to recover a missing or outdated metamodel [30]. The process involves converting the instance models in XML format to a domain-specific language, then performing existing grammar inference techniques on those DSL programs. The authors use their previously developed evolutionary approach [57] to do the actual inference, then recreate a metamodel in XML format from the result so the recovered metamodel can be loaded into a modeling tool. Liu et al. have recently extended this system to handle models with a more complex and segmented organizational structure [44]. The authors refer to these as multi-tiered domains because they support multiple viewpoints into the model.

Two similar problems are grammar convergence [37] and grammar recovery [36], both which involve finding grammars for a variety of software artifacts. The goal of grammar convergence is to establish and maintain a mapping between software artifact structures in different formats that embody the same domain knowledge. Grammatical inference can aid in the early steps of this process to produce a grammar for each knowledge representation by examining available concrete examples. Existing grammar transformation and convergence techniques can then be used on the resulting source grammars to establish a unified grammar.

Grammar recovery can be viewed as a more general version of the grammar inference problem because it seeks to recover a grammar from sources such as compilers, reference manuals and written specifications, in addition to concrete program examples. The effort by Lämmel and Verhoef to recover a VS COBOL II grammar includes leveraging visual syntax diagrams from the manual [36]. These diagrams give clues about the shape of the target grammar’s derivation tree, knowledge that is known to greatly improve the accuracy of grammatical inference techniques. For example, reversible context-free and non-counting context-free languages are known to be identifiable from positive examples with these types of structured strings [52, 16]. Furthermore, structured strings can be used to identify any context-free language in polynomial time with a membership and equivalence oracle [51]. In the case of grammar recovery, an existing compiler for the language (even without the compiler source code) may be used as a membership oracle.

5 Related Surveys

Many surveys of grammatical inference have been written to introduce newcomers to the field and summarize the state of the art. Most give a thorough overview of grammatical inference in general, but each emphasise different aspects of the literature.

Fu and Booth (1986) give a detailed technical description of some early inference algorithms and heuristics with an emphasis on pattern recognition examples to demonstrate its relevance [22]. This survey is heavy on technical definitions and grammatical notation, suitable for someone with prior knowledge in formal languages who prefers to get right into the algorithms and techniques of grammatical inference.

Vidal (1994) provides a concise but thorough overview of the learning models and language classes in grammatical inference, with ample citations for follow-up investigation [58]. He presents each learning model in the context of fundamental learnability results in the field as well as their practical applications without getting too deeply into the details of each learning model.

Dana Ron's doctoral thesis (1995) on the learning of deterministic and probabilistic finite automata primarily investigates PAC learning as it relates to identifying DFAs, and describes practical applications of its use [50]. Although not exhaustive of grammatical inference in general, this thesis is a good reference for someone specifically interested in DFA inference.

Lee (1996) presents an extensive survey on the learning of context-free languages, including those that have non-grammar formalisms [42]. She discusses approaches that learn from both text and structured data, making it relevant to software engineering induction problems.

Sakakibara (1997) provides an excellent overview of the field with an emphasis on computational complexity, learnability, and decidability [53]. He covers a wide range of grammar classes including DFAs, context-free grammars and their probabilistic counterparts. This survey is roughly organized by the types of language classes being learned.

Colin de la Higuera (2000) gives a high-level and approachable commentary on grammatical inference including its historical progress and achievements [18]. He highlights key issues and unsolved problems, and describes some promising avenues of future research. This commentary is not meant as a technical introduction to inference techniques nor an exhaustive survey, and therefore contains no mathematical or formal notation. It rather serves as a quick and motivational read for anyone interested in learning about grammatical inference. A similar piece on grammatical inference is written by Honavar and de la Higuera (2001) for a special issue of the *Machine Learning* journal (volume 44), emphasizing the cross-disciplinary nature of the field.

Cicchello and Kremer (2003) and Bugalho and Oliveira (2005) survey DFA inference algorithms in depth, with excellent explanations about augmented prefix tree acceptors, state merging, the red-blue framework, search heuristics, and performance comparisons of state of the art DFA inference algorithms [13, 10].

de la Higuera (2005) provides an excellent guide to grammatical inference, geared toward people (not necessarily experts in formal languages or computational linguistics) who think grammatical inference may help them solve their particular problem [19]. He gives a general roadmap of the field, examples of how grammatical inference has been used in existing applications, and provides many useful references for further investigation by the reader.

Pieter Adriaans and Menno van Zaanen (2006) compare grammatical inference from three different perspectives: linguistic, empirical, and formal [1]. They introduce the common learning models and broad learnability results in the framework of each perspective, and comment on how these perspectives overlap. This survey is useful for someone who comes from a linguistic, empirical, or formal languages background and wishes to learn about grammatical inference.

6 Future Direction and Challenges

Software engineers are finding a variety of uses for grammatical inference in their work, but grammatical inference is still relatively rare in the field. The theoretical work in grammatical inference is largely disconnected from these practical uses because implementers tend to use domain specific knowledge to craft custom solutions. Such solutions, while successful in some cases, usually ignore the powerful algorithms developed by the theoretical GI community. Domain knowledge should continue to be exploited – we are not advocating otherwise – but domain knowledge needs to be translated into a form general-purpose inference algorithms can use. We believe this is the biggest challenge currently facing software engineers wanting to use grammatical inference in their applications: how to map their domain knowledge to theoretical GI constraints.

Constraints can be imposed in several ways, such as simplifying the grammar class to learn, providing negative samples, adding a membership and/or equivalence oracle where none existed, or partially structuring the input data. Often these constraints are equivalent to some existing structural knowledge or implicit assumptions about the input data, but identifying these equivalences is nontrivial.

```

while limit > a do
  begin
    if a > max then max = a;
    a := a + 1
  end

```

Fig. 2. A sample program in an unknown Pascal-like language

We motivate this approach with a concrete example, inspired by an example from Sakakibara [52]. Suppose you have a collection of computer programs written in an unknown language with a Pascal-like syntax and wish to infer a

grammar from the collection. For clarity, Figure 2 shows a sample program in the collection and Figure 3 shows the target grammar to learn (a subset of the full Pascal grammar).

$$\begin{aligned}
 \textit{Statement} &\rightarrow \textit{Ident} := \textit{Expression} \\
 \textit{Statement} &\rightarrow \mathbf{while} \textit{Condition} \mathbf{do} \textit{Statement} \\
 \textit{Statement} &\rightarrow \mathbf{if} \textit{Condition} \mathbf{then} \textit{Statement} \\
 \textit{Statement} &\rightarrow \mathbf{begin} \textit{Statementlist} \mathbf{end} \\
 \textit{Statementlist} &\rightarrow \textit{Statement}; \textit{Statementlist} \\
 \textit{Statementlist} &\rightarrow \textit{Statement} \\
 \textit{Condition} &\rightarrow \textit{Expression} > \textit{Expression} \\
 \textit{Expression} &\rightarrow \textit{Term} + \textit{Expression} \\
 \textit{Expression} &\rightarrow \textit{Term} \\
 \textit{Term} &\rightarrow \textit{Factor} \\
 \textit{Term} &\rightarrow \textit{Factor} \times \textit{Term} \\
 \textit{Factor} &\rightarrow \textit{Ident} \\
 \textit{Factor} &\rightarrow (\textit{Expression})
 \end{aligned}$$

Fig. 3. The target grammar for the Pascal-like language [52]

At first glance this inference problem seems too difficult to solve. It is a context-free grammar with positive samples only, and Gold proved learning a superfinite language in the limit from positive-only samples is impossible [24]. Even with the addition of negative samples there is no known algorithm to efficiently learn a context-free language.

On closer inspection, however, there is additional structural information in the input samples, hidden in a place grammarware authors and parsers are trained to ignore – the whitespace. By taking into account line breaks and indented sections of source code in the input samples, a structured string can be constructed for each program. If we further assume the target grammar is reversible then we can apply a result by Sakakibara, who showed reversible context-free grammars can be learned in polynomial time from positive structured strings [52].

This particular solution depends on two assumptions: (1) all the input samples have meaningful and consistent whitespace formatting, and (2) the target grammar is in the class of reversible context-free grammars. The assumption that the target grammar is reversible context-free is reasonable, as many DSLs would fit this criterion. The Pascal subset grammar in Figure 3 is in fact reversible context-free, but full Pascal is not because adding a production rule like $\textit{Factor} \rightarrow \textit{Number}$ to this grammar violates the criteria of reversibility [52].

Leveraging domain knowledge and structural assumptions is quite powerful when inferring grammars from examples and should be encouraged, but at

present mapping this domain-specific knowledge to abstract constructs in grammatical inference research requires some creativity and awareness of theoretical results in the field. Allowing the extensive work done in the theoretical grammatical inference community to bear on specific applications of GI would be a great boon to software engineering.

7 Conclusion

In grammatical inference, an inference algorithm must find and use common patterns in example sentences to concisely represent an unknown language in the form of a grammar. This process spans two axes of complexity: the language class to be learned and the learning model employed.

The theoretical foundations of grammatical inference are now well established thanks to contributions by Gold, Angluin and others. The state of the art, however, still has plenty of room to grow, and Colin de la Higuera identifies ten open problems on the theoretical side of grammatical inference that he believes are important to solve going forward [20].

In practice, assumptions can often be made which are not possible in a purely theoretical setting because a specific problem domain has limited scope, allowing for a better outcome than one would expect from simply applying the smallest enclosing theoretical result. Some work has already been done to investigate this relationship deeper, such as that by Kermorvant and de la Higuera [34] and Cano et al [12]. It would be valuable to find a widely applicable technique to equate domain assumptions with either a restriction in the class of language to learn, or an augmentation of the learning model.

Theoretical grammatical inference research continues to advance in many different directions: the language classes being learned, the learning models in use, the criteria for a successful inference, and the efficiency of the inference algorithms. Existing applications for grammatical inference are continually refined and new applications are found in a wide variety of disciplines.

References

1. Adriaans, P., van Zaanen, M.: Computational grammatical inference. *Studies in Fuzziness and Soft Computing* **194** (2006) 187–203
2. Angluin, D.: Inductive inference of formal languages from positive data. *Information and Control* **45**(2) (1980) 117–135
3. Angluin, D.: A note on the number of queries needed to identify regular languages. *Information and Control* **51**(1) (1981) 76–87
4. Angluin, D.: Inference of reversible languages. *Journal of the ACM (JACM)* **29** (1982) 741–765
5. Angluin, D.: Learning regular sets from queries and counterexamples. *Information and Computation* **75** (November 1987) 87–106
6. Angluin, D.: Queries and concept learning. *Machine Learning* **2**(4) (1988) 319–342
7. Angluin, D.: Negative results for equivalence queries. *Machine Learning* **5**(2) (July 1990) 121–150

8. Angluin, D., Kharitonov, M.: When won't membership queries help? In: Proceedings of the twenty-third annual ACM symposium on Theory of computing. STOC '91, New York, NY, USA, ACM (1991) 444–454
9. Ates, K., Kukluk, J., Holder, L., Cook, D., Zhang, K.: Graph grammar induction on structural data for visual programming. In: 18th IEEE International Conference on Tools with Artificial Intelligence, 2006. ICTAI '06. (November 2006) 232–242
10. Bugalho, M., Oliveira, A.L.: Inference of regular languages using state merging algorithms with search. *Pattern Recogn.* **38**(9) (September 2005) 1457–1467
11. Burago, A.: Learning structurally reversible context-free grammars from queries and counterexamples in polynomial time. In: Proceedings of the seventh annual conference on Computational learning theory. COLT '94, New York, NY, USA, ACM (1994) 140–146
12. Cano, A., Ruiz, J., García, P.: Inferring subclasses of regular languages faster using RPNI and forbidden configurations. In: Proceedings of the 6th International Colloquium on Grammatical Inference: Algorithms and Applications. ICGI '02, London, UK, UK, Springer-Verlag (2002) 28–36
13. Cicchello, O., Kremer, S.C.: Inducing grammars from sparse data sets: a survey of algorithms and results. *J. Mach. Learn. Res.* **4** (December 2003) 603–632
14. Clark, A.: Distributional learning of some context-free languages with a minimally adequate teacher. In: Proceedings of the 10th international colloquium conference on Grammatical inference: theoretical results and applications. ICGI'10, Berlin, Heidelberg, Springer-Verlag (2010) 24–37
15. Clark, A., Eyraud, R., Habrard, A.: A polynomial algorithm for the inference of context free languages. In: Proceedings of the 9th international colloquium on Grammatical Inference: Algorithms and Applications. ICGI '08, Berlin, Heidelberg, Springer-Verlag (2008) 29–42
16. Crespi-Reghezzi, S., Guida, G., Mandrioli, D.: Noncounting context-free languages. *Journal of the ACM (JACM)* **25**(4) (October 1978) 571–580
17. Crespi-Reghezzi, S., Melkanoff, M.A., Lichten, L.: The use of grammatical inference for designing programming languages. *Communications of the ACM* **16** (1973) 83–90
18. de la Higuera, C.: Current trends in grammatical inference. In: Proceedings of the Joint IAPR International Workshops on Advances in Pattern Recognition, London, UK, Springer-Verlag (2000) 28–31
19. de la Higuera, C.: A bibliographical study of grammatical inference. *Pattern Recognition* **38** (September 2005) 1332–1348
20. de la Higuera, C.: Ten open problems in grammatical inference. In Sakakibara, Y., Kobayashi, S., Sato, K., Nishino, T., Tomita, E., eds.: *Grammatical Inference: Algorithms and Applications*. Volume 4201 of *Lecture Notes in Computer Science*. Springer Berlin / Heidelberg (2006) 32–44
21. Dubey, A., Jalote, P., Aggarwal, S.: Learning context-free grammar rules from a set of programs. *Software, IET* **2**(3) (2008) 223–240
22. Fu, K.S., Booth, T.L.: Grammatical inference: introduction and survey\part i. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **8** (May 1986) 343–359
23. Fürst, L., Mernik, M., Mahnic, V.: Graph grammar induction as a parser-controlled heuristic search process, Budapest, Hungary (October 2011)
24. Gold, E.M.: Language identification in the limit. *Information and Control* **10**(5) (1967) 447–474
25. Gold, E.M.: Complexity of automaton identification from given data. *Information and Control* **37**(3) (1978) 302–320

26. Hrnčič, D., Mernik, M., Bryant, B.R.: EMBEDDING DSLS INTO GPLS: a GRAMMATICAL INFERENCE APPROACH *. *Information Technology And Control* **40**(4) (December 2011)
27. Hrnčič, D., Mernik, M., Bryant, B.R., Javed, F.: A memetic grammar inference algorithm for language learning. *Applied Soft Computing* **12**(3) (March 2012) 1006–1020
28. Ishizaka, H.: Polynomial time learnability of simple deterministic languages. *Machine Learning* **5**(2) (July 1990) 151–164
29. Javed, F., Mernik, M., Bryant, B.R., Gray, J.: A grammar-based approach to class diagram validation. (2005)
30. Javed, F., Mernik, M., Gray, J., Bryant, B.R.: MARS: a metamodel recovery system using grammar inference. *Inf. Softw. Technol.* **50**(9-10) (August 2008) 948–968
31. Javed, F., Mernik, M., Sprague, A., Bryant, B.: Incrementally inferring context-free grammars for domain-specific languages. *Proceedings of the Eighteenth International Conference on Software Engineering and Knowledge Engineering (SEKE'06)* (2006) 363–368
32. Juillé, H., Pollack, J.B.: A stochastic search approach to grammar induction. In: *Proceedings of the 4th International Colloquium on Grammatical Inference. ICGI '98*, London, UK, UK, Springer-Verlag (1998) 126–137
33. Kearns, M., Li, M., Pitt, L., Valiant, L.: On the learnability of boolean formulae. In: *Proceedings of the nineteenth annual ACM symposium on Theory of computing. STOC '87*, New York, NY, USA, ACM (1987) 285–295
34. Kermorvant, C., Higuera, C.D.L.: Learning languages with help. In: *Proceedings of the 6th International Colloquium on Grammatical Inference: Algorithms and Applications. ICGI '02*, London, UK, UK, Springer-Verlag (2002) 161–173
35. Kong, J., Ates, K., Zhang, K., Gu, Y.: Adaptive mobile interfaces through grammar induction. In: *20th IEEE International Conference on Tools with Artificial Intelligence, 2008. ICTAI '08. Volume 1.* (November 2008) 133–140
36. Lämmel, R., Verhoef, C.: Semi-automatic grammar recovery. *Softw. Pract. Exper.* **31**(15) (December 2001) 1395–1448
37. Lämmel, R., Zaytsev, V.: An introduction to grammar convergence. In: *Proceedings of the 7th International Conference on Integrated Formal Methods. IFM '09*, Berlin, Heidelberg, Springer-Verlag (2009) 246–260
38. Lang, K.J.: Faster algorithms for finding minimal consistent DFAs. *Technical report* (1999)
39. Lang, K.J., Pearlmutter, B.A., Price, R.A.: Results of the abbadingo one DFA learning competition and a new evidence-driven state merging algorithm. In: *Proceedings of the 4th International Colloquium on Grammatical Inference*, London, UK, Springer-Verlag (1998) 1–12
40. Langley, P., Stromsten, S.: Learning context-free grammars with a simplicity bias. *Proceedings of the Eleventh European Conference on Machine Learning* (2000) 220–228
41. Larus, J.R.: Whole program paths. In: *ACM SIGPLAN Notices. PLDI '99*, New York, NY, USA, ACM (1999) 259–269
42. Lee, L.: Learning of context-free languages: A survey of the literature. *REP* (1996) 12–96
43. Li, M., Vitányi, P.M.B.: Learning simple concepts under simple distributions. *Siam Journal of Computing* **20** (1991) 911–935
44. Liu, Q., Bryant, B.R., Mernik, M.: Metamodel recovery from multi-tiered domains using extended MARS. In: *Proceedings of the 2010 IEEE 34th Annual Computer*

- Software and Applications Conference. COMPSAC '10, Washington, DC, USA, IEEE Computer Society (2010) 279–288
45. Memon, A.U.: Log File Categorization and Anomaly Analysis Using Grammar Inference. Master of science, Queen's University (2008)
 46. Mernik, M., Hrnčić, D., Bryant, B., Sprague, A., Gray, J., Liu, Q., Javed, F.: Grammar inference algorithms and applications in software engineering. In: Information, Communication and Automation Technologies, 2009. ICAT 2009. XXII International Symposium on. (October 2009) 1–7
 47. Nevill-Manning, C.G., Witten, I.H.: Identifying hierarchical structure in sequences: a linear-time algorithm. *Journal of Artificial Intelligence Research* **7**(1) (September 1997) 67–82
 48. Oncina, J., García, P.: Identifying regular languages in polynomial time. In: Advances in Structural and Syntactic Pattern Recognition - Proceedings of the International Workshop on Structural and Syntactic Pattern Recognition, Bern, Switzerland (1992) 99–108
 49. Pitt, L., Valiant, L.G.: Computational limitations on learning from examples. *Journal of the ACM (JACM)* **35**(4) (October 1988) 965–984
 50. Ron, D.: Automata Learning and its Applications. PhD thesis, Hebrew University (1995)
 51. Sakakibara, Y.: Learning context-free grammars from structural data in polynomial time. *Theoretical Computer Science* **76**(2-3) (November 1990) 223–242
 52. Sakakibara, Y.: Efficient learning of context-free grammars from positive structural examples. *Information and Computation* **97**(1) (1992) 23–60
 53. Sakakibara, Y.: Recent advances of grammatical inference. *Theoretical Computer Science* **185** (October 1997) 15–45
 54. Trakhtenbrot, B.A., Barzdin, Y.M.: Finite Automata: Behaviour and Synthesis. North-Holland Publishing Company, Amsterdam (June 1973)
 55. Valiant, L.G.: A theory of the learnable. *Communications of the ACM* **27** (1984) 1134–1142
 56. Črepinšek, M., Mernik, M., Bryant, B.R., Javed, F., Sprague, A.: Inferring context-free grammars for domain-specific languages. *Electronic Notes in Theoretical Computer Science* **141**(4) (December 2005) 99–116
 57. Črepinšek, M., Mernik, M., Javed, F., Bryant, B.R., Sprague, A.: Extracting grammar from programs: evolutionary approach. *ACM SIGPLAN Notices* **40** (2005) 39–46
 58. Vidal, E.: Grammatical inference: An introductory survey. In Carrasco, R., Oncina, J., eds.: *Grammatical Inference and Applications*. Volume 862 of *Lecture Notes in Computer Science*. Springer Berlin / Heidelberg (1994) 1–4
 59. Yokomori, T.: Polynomial-time learning of very simple grammars from positive data. In: *Proceedings of the fourth annual workshop on Computational learning theory*, San Francisco, CA, USA, Morgan Kaufmann Publishers Inc. (1991) 213–227
 60. Yokomori, T.: On polynomial-time learnability in the limit of strictly deterministic automata. *Machine Learning* **19**(2) (1995) 153–179