



Probabilistic Pattern Matching and the Evolution of Stochastic Regular Expressions

BRIAN J. ROSS

Department of Computer Science, Brock University, St. Catharines, Ontario, Canada L2S 3A1

bross@cosc.brocku.ca

Abstract. The use of genetic programming for probabilistic pattern matching is investigated. A stochastic regular expression language is used. The language features a statistically sound semantics, as well as a syntax that promotes efficient manipulation by genetic programming operators. An algorithm for efficient string recognition based on approaches in conventional regular language recognition is used. When attempting to recognize a particular test string, the recognition algorithm computes the probabilities of generating that string and all its prefixes with the given stochastic regular expression. To promote efficiency, intermediate computed probabilities that exceed a given cut-off value will pre-empt particular interpretation paths, and hence prune unconstructive interpretation. A few experiments in recognizing stochastic regular languages are discussed. Application of the technology in bioinformatics is in progress.

Keywords: stochastic regular expressions, genetic programming

1. Introduction

Language inference is a classical problem in machine learning, and continues to be an important and active research topic. The basic problem is, given a set of example behaviours or strings, automatically infer a corresponding language (grammar, automata, expression, . . .) which generates or recognizes those examples. Genetic algorithms (GA) and genetic programming (GP) have been applied towards language inference, with varying degrees of success. Although successful inference is possible, the generic inference problem is not entirely well-suited for solution by evolutionary search. There are a number of reasons for this. For example, some genome encodings do not preserve useful language characteristics during crossover. Even small local changes to such genomes can be catastrophic, which does not lend itself well to genetic reproduction and evolutionary search.

An even more acute weakness is that “all or nothing” problems such as the language inference problem are not entirely natural for GP. An acceptable language inference minimally requires that the solution language

correctly recognize all positive test cases, and reject all negative ones. This essential criteria may also be supplemented by efficiency concerns, such as a relatively small number of states or grammar rules. The resulting search space is a difficult one to navigate with evolutionary techniques, due to these stringent requirements for language correctness and completeness. On the other hand, it is generally recognized in the GP community that problems which require an “acceptably close” solution are typically the best candidates for successful solution with GP. Pragmatically speaking, giving the fitness function a larger degree of freedom for evaluating a successful solution will substantially increase the chances of the discovery of acceptable solutions.

This research addresses the inference of stochastic regular languages using genetic programming. Stochastic languages are formal languages with probability distribution associated with the language set. The stochastic language inference problem is similar to the classical inference problem, with the additional requirement that the distribution of strings recognized by the stochastic language conform to some desired

target distribution. At first, this may seem intuitively more complex than non-stochastic language inference, since it is unclear what impact the determination of probability distributions has on the inference. It turns out, however, that the inclusion of string distributions can simplify the inference problem. Hypothesized languages are now allowed to generate erroneous strings so long as they fall within an acceptably small probability of occurrence. In other words, the use of language distributions introduces a more generous degree of freedom for generated solutions. This is ideal in a GP setting, as it simplifies the search space substantially for evolutionary search.

The target language used here is a probabilistic regular expression language, henceforth called Stochastic Regular Expressions (or SRE). Although theoretically weaker than stochastic context-free languages studied elsewhere, it was nevertheless chosen due to both its amenability to concise GP representation, and its ability to naturally solve the substantial number of problems in the “regular language domain”. The stochastic regular expression language is closely related to stochastic regular grammars and stochastic finite automata, the latter commonly referred to as Hidden Markov Models in the literature.

Some SRE language implementation issues had to be addressed before GP could be successfully applied to stochastic regular expression problems. Firstly, an efficient implementation of SRE interpretation was necessary. Interpretation of an SRE expression requires that the probability of recognizing a given string is generated. Since intermediate probabilities would be computed during the interpretation of a string, these values can be used to terminate or prune unproductive interpretation paths whose probabilities are smaller than some supplied cut-off probability. Given the extensive testing that is necessary during fitness evaluation, such pruning greatly increases the speed of GP runs. The SRE language is implemented in a grammatical GP system, which permitted the use of syntactic language constraints to further enhance evolution efficiency.

Two example experiments proved that probabilistic language inference is indeed possible with SRE and GP. The more complex of these experiments indicated that the complex search space often resulted in premature convergence. A minor language enhancement to this experiment resulted in failed inferences by the GP system. From this experience, it can be deduced that the fitness evaluation strategy used here is not a general purpose solution to all stochastic language problems,

but rather, is suitable to a class of stochastic regular languages whose members are structurally related to one another.

An outline of the paper is as follows. Related work is reviewed in Section 2. Section 3 defines the syntax and semantics of the stochastic regular expression language, and discusses the algorithm for processing SRE expressions. Section 4 outlines the genetic programming system used. Two example experiments are discussed in Section 5. A discussion and future directions conclude the paper in Section 6.

2. Related Work

Formal language induction has a long history as a fundamental problem in machine learning [1–4]. The specialized topic of stochastic languages has also been studied for some time [5]. A stochastic grammar differs from a conventional grammar in that each grammar rule is marked with a probability associated with its use, and the set of probabilities for a grammar encode a probability distribution for the resulting derived language. [6] has an extensive treatment of stochastic grammars, their derivation, and their application in pattern recognition. Stochastic grammars are also more complex than their non-stochastic kin, as the distributions inherent with the language introduce a new dimension of membership criteria. For example, all context free languages are also stochastic context free languages (all probabilities are 1); however, there may be many stochastic context free languages having essentially the same membership set, but vastly different distributions over that set. Language equivalence issues are therefore more discriminating than in a non-stochastic setting. Stochastic context free languages enjoy both expressivity and tractable properties, for example, the existence of useful inference algorithms [7]. They have also found practical use in language processing [8].

Stochastic regular languages, albeit descriptively weaker than stochastic context-free languages, have also found their practical niche in applications. Regular languages are definable by finite automata, regular grammars, and regular expressions [9]. Similarly, stochastic regular languages are defined by stochastic versions of these three representations. Examples of work in stochastic grammar inference are in [10–13]. Stochastic finite automata are defined in terms of Hidden Markov Models (HMM) [14]. An HMM is a finite automaton with probabilities marking the transition links between nodes. Each node is connected to all

other nodes, and so the network itself is maximally connected. When particular transitions are not required, the probabilities associated with those nodes are set to zero. HMMs have found extensive use in language and speech processing [8, 15]. Strangely enough, stochastic regular expressions have not been extensively studied; one example paper is [16].

Language inference has been successfully done using genetic algorithms (GA) and genetic programming (GP). The distinguishing difference between GA and GP approaches is one of denotation: a pure GA uses a binary encoding for the genome, while a GP uses a variable-sized parse tree. Some of the following use encodings with characteristics of both approaches.

With respect to regular languages, an early work in evolving finite automata is by Zhou and Grefenstette [17]. They used a GA with a binary encoding of the automata as a set of state transitions, capped at a size of 8 states. A weakness of this encoding is that the represented automata are susceptible to destructive effects during crossover and mutation. Their unspecified fitness function scores language performance (ability to accept positive strings and reject negative examples) and automata size.

Dunay et al.'s approach [18] is similar to [17], except that finite automata are denoted in GP-style nested S-expression notation.

Dupont [19] proposes an automata-theoretic partition representation for regular languages. This has the advantage of preserving language properties of chromosomes during GA reproduction, unlike the more fragile FA representation in [17]. His fitness function scores both language performance and automata size. He successfully evolved a large set of regular languages, including the benchmark Tomita languages [20].

Brave [21] uses an abstract "cellular encoding" representation for deterministic FA's, which builds the network structure of a FA during interpretation. The intention of this denotation is to preserve structural properties of a language during evolutionary reproduction. His automata are embellished with boolean operators which permit automata composition. The fitness function tallied the number of correctly classified sentences. All but one of the Tomita languages were successfully inferred using this technique.

Longshaw [22] adopts a straight-forward state-transition representation for automata. However, his GA uses a population seeded with correct but overly general automata. Specialized reproduction operators

manipulate automata by duplicating or refining states. The overall intention is to refine the general automata into a more specific one for the language in question. His fitness function scores example classification performance and automata size.

Svingen [23] applies GP on regular expressions. Regular expressions are directly encoded as program trees, and fitness is based on correct example classification. He successfully evolved the Tomita languages.

Context-free languages have also been studied. Wyard [24] uses a GA to evolve context-free grammars. Chromosomes takes the form of lists of production rules, which guarantees correctness at all times. The fitness function scores example classification performance. Two simple CFG's were successfully evolved.

Lankhorst [25] applies a vector encoding to represent grammar productions. His fitness function is more involved than most others, as it scores example classification performance, the length of substrings of examples correctly classified, the degree of determinism of grammars, and the ability of the grammar to generate correct strings not included in the example set. These additional evaluation considerations give the GA more information with which to drive evolution. He applied the GA to a number of CFG and regular languages.

Lucas [26] suggests a binary-encoded normal form for CFG productions, which preserves language properties during reproduction, and promotes convergence. His fitness strategy scores example classification and grammar size.

Sen and Janakiraman [27] apply a GA towards inferring deterministic pushdown automata, which is an alternative to the grammar representation for CFG's. Fitness scores example recognition performance, and whether the PDA attempts to erroneously 'pop' an empty stack. Lankhorst [28] extends this idea towards nondeterministic pushdown automata. His fitness additionally considers prefix sizes and the stack size after a string has been consumed.

Dunay and Petry [29] utilize a Turing machine representation in their GA experiments. Although this powerful notation can denote the entire set of languages in the Chomsky hierarchy, it does not necessary mean that search will be easy to accomplish, given the inherent enormity of the search space in question. To solve some relatively simple examples of regular, context-free and context-sensitive languages, they used a compositional approach, in which the GA had access to TM building blocks evolved in earlier runs.

Finally, the evolution of stochastic languages has been studied. Schwehm and Ost [30] use a GA for evolving stochastic regular languages. Two different encodings are studied—production rules with probabilities, and quotient automata. The fitness function uses grammar complexity (number of productions), a modified χ^2 test for distribution conformance, and a measure of the grammar’s ability to accept prefixes of the target grammar. A few experiments were performed, and their GA performance compares well with standard regular-language inference algorithms.

Kammeyer and Belew [31] investigate the application of GA’s to evolve stochastic context-free grammars. They use a liberal representation for grammars in which correct grammars are parsed from the genome when evaluated; this permits intron or junk material to be included in chromosomes. The fitness function evaluates the size of test example prefixes consumed by a grammar, and uses cross-entropy to evaluate distribution conformance. They also use a local search technique for finding production probabilities during evolution. A couple of CFG’s were successfully evolved.

3. Stochastic Regular Expressions

3.1. Language Definition

The target language for the GP system is stochastic regular expressions, or SRE. The language is very similar to one in [16], which is used for modeling the qualitative behaviour of stochastic discrete event systems. Amongst other properties, they prove that probabilistic regular language operations such as choice, concatenation, and Kleene-closure forms a closed language, and hence an algebra. Although a few basic properties will be illustrated here, the reader is referred to [16] for further details. It is assumed the reader is familiar with basic concepts from formal language theory [9].

Two language variations, SRE and Guarded SRE (or gSRE), are used. We first define SRE. Let α range over alphabet Σ , E range over SRE expressions, n range over positive integers ($0 \leq n \leq 1000$), and f range over decimal values with a precision of 2 decimal places ($0 \leq f < 1.00$). The syntax of SRE is recursively defined as:

$$E ::= \alpha \mid \sum_i E_i(n_i) \mid E_1 : E_2 \mid E^{*f} \mid E^{+f}$$

Without loss of generality, the empty string ϵ is not included in the alphabet.

The operators have the following meaning:

1. *Atomic action* α : The action α is generated.
2. *Choice* $\sum_i E_i(n_i)$: This denotes a probabilistic choice of terms. Each choice expression E_i can be chosen with a probability:

$$\frac{n_i}{\sum_j n_j}$$

For example, given the expression $E_1(3) + E_2(5)$, the term E_1 can be chosen with a probability of 3/8 and E_2 with a probability of 5/8.

3. *Concatenation* “ $E_1 : E_2$ ”: Term E_1 is interpreted, followed by that of E_2 .
4. *Kleene Closure* E^{*f} : Term E can be repeatedly executed 0 or more times, and each iteration occurs with a probability of f . The probability of E terminating execution is $1 - f$.
5. *+Closure* E^{+f} : Term E executes once, after which it repeatedly executes 0 or more times using the same probability scheme as Kleene closure. +Closure is an abbreviation for the following:

$$E^{+f} \equiv E : E^{*f}$$

The Guarded SRE language is identical to SRE, except that a guarded choice operator is used instead of the general choice in 2 above:

6. *Guarded Choice* $\sum_i E'_i(n_i)$, where $E' = (\alpha_i : E_i)$ or $E' = \alpha_i$, and $\forall \alpha_i, \alpha_j : \alpha_i \neq \alpha_j$: Here, each term in the choice expression is either prefixed with a unique atomic action that is found nowhere else in the expression, or consists of a unique action by itself. This makes guarded choice deterministic, unlike SRE’s nondeterministic choice.

Note that, even with guarded choice, gSRE is still a nondeterministic language, since the closure operators are nondeterministic. The rest of the discussion in this section pertains equally to both SRE and gSRE.

A derivation of a conventional regular expression E is the set of sentences, or strings over the alphabet, derivable from it. This defines the language $L(E)$ of E . This notion of language derivation is similarly applicable to SRE, except that each string has a probability value associated with it, and hence the language itself is associated with a probability distribution of its members.

Alternatively, an intuitive way to consider SRE expressions is that every expression defines a specific

probability function over strings in Σ^* :

$$E : \Sigma^* \rightarrow p \quad (0 \leq p \leq 1)$$

Using a denotational semantics style of representation [32], the probability function for SRE expression E is denoted by $\llbracket E \rrbracket$, and its application to a particular string s is denoted $\llbracket E \rrbracket s$, which denotes the probability associated with string s in the language $L(E)$.

A probability function model of SRE is now given. Let $s = \alpha_1, \dots, \alpha_n \in \Sigma^*$.

- Atomic actions:

$$\begin{aligned} \llbracket \alpha \rrbracket \beta &= 1 & \text{if } \alpha = \beta \\ \llbracket \alpha \rrbracket \beta &= 0 & \text{if } \alpha \neq \beta \end{aligned} \quad (1)$$

- Choice (including guarded choice):

$$\left[\sum_i E_i(n_i) \right] s = \sum_k \left(\frac{n_k}{\sum_j n_j} \cdot \llbracket E_k \rrbracket s \right) \quad (2)$$

Since every term might recognize s , the overall probability for a choice expression is the sum of all the term probabilities with respect to s .

- Concatenation:

$$\begin{aligned} \llbracket E_1 : E_2 \rrbracket s &= \sum_{i=1}^n (\llbracket E_1 \rrbracket \alpha_1 \dots \alpha_i \cdot \llbracket E_2 \rrbracket \alpha_{i+1} \dots \alpha_n) \\ &+ \llbracket E_1 \rrbracket s \cdot \llbracket E_2 \rrbracket \epsilon \\ &+ \llbracket E_1 \rrbracket \epsilon \cdot \llbracket E_2 \rrbracket s \end{aligned} \quad (3)$$

In the first summation, s is decomposed into two substrings, each of which may be consumed by a concatenated expression. Even though one term may recognize its substring argument, if the other term does not recognize its respective substring, then that term returns a probability of 0, and the overall probability for that instance of decomposition is 0. The rest of the formula represents the cases when one entire expression consumes s , while the other consumes ϵ . If these other cases do not succeed, then they return 0.

- Kleene closure:

$$\begin{aligned} \llbracket E^{*f} \rrbracket \epsilon &= 1 - f \\ \llbracket E^{*f} \rrbracket s &= \sum_{i=1}^n (f \cdot \llbracket E \rrbracket \alpha_1 \dots \alpha_{i-1} \cdot \llbracket E^{*f} \rrbracket \alpha_i \dots \alpha_n) \\ &+ f \cdot \llbracket E \rrbracket s \cdot \llbracket E^{*f} \rrbracket \epsilon \quad s \neq \epsilon \end{aligned} \quad (4)$$

The first formula accounts for empty strings, as the only way an iterated expression should recognize an empty string is by not iterating. The other formula recursively defines the general case. Here, one iteration of E will consume some portion of s , and the rest of s is consumed by further iterations. The final term in this formula represents when the first iteration consumes the entire string. It is assumed that an iteration of a loop always consumes some non-empty string. Otherwise, the semantic model would have to account for Kleene closure iterating indefinitely on an argument, which is not useful behaviour.

- +Closure:

$$\begin{aligned} \llbracket E^{+f} \rrbracket s &= \sum_{i=1}^n (\llbracket E \rrbracket \alpha_1 \dots \alpha_i \cdot \llbracket E^{*f} \rrbracket \alpha_{i+1} \dots \alpha_n) \\ &+ \llbracket E \rrbracket s \cdot \llbracket E^{*f} \rrbracket \epsilon \end{aligned} \quad (5)$$

This is similar to the non-empty argument formula for Kleene closure, except that the expression E will consume part of the string before iterations commence. This can be seen by the lack of f value in the formula.

The nondeterministic nature of regular expressions is modeled in the above by multiple argument decomposition in the concatenation and closure operators. Nondeterminism can also arise in the (nonguarded) choice operator.

Membership in SRE is reflected by SRE expressions returning non-zero probabilities for particular strings:

$$\begin{aligned} s \in L(E) &\text{ iff } \llbracket E \rrbracket s > 0 \\ s \notin L(E) &\text{ iff } \llbracket E \rrbracket s = 0 \end{aligned}$$

Definition 3.1. All probability functions pf must adhere to the following two characteristics [33]:

- (i) for all x_i in the sample space of the experiment:

$$\sum_i pf(x_i) = 1 \quad (6)$$

- (ii) For every event x_i :

$$0 \leq pf(x_i) \leq 1 \quad (7)$$

Consequently, if SRE expressions are to define well-formed probability functions, all expressions must similarly respect these requirements.

Theorem 3.1. *The SRE operators are well-formed probability functions.*

Proof: The proof uses structural induction on SRE expressions. We show conditions (i) and (ii) of Definition 3.1 hold for all operators. Let $s \in \Sigma^*$.

- (a) *Atomic actions:* trivially.
 (b) *Choice:* (i) From Eqs. (2) and (6):

$$\sum_i \sum_k \left(\frac{n_k}{\sum_j n_j} \cdot \llbracket E_k \rrbracket_{s_i} \right)$$

By the induction hypothesis,

$$\sum_i \llbracket E_k \rrbracket_{s_i} = 1.$$

Thus we have,

$$\sum_k \frac{n_k}{\sum_j n_j} = 1.$$

- (ii) From Eq. (2), the greatest value for the sum

$$\sum_k \left(\frac{n_k}{\sum_j n_j} \cdot \llbracket E_k \rrbracket_s \right)$$

occurs when $\llbracket E_k \rrbracket_s = 1$ for all k . In this case, the sum reduces to

$$\frac{\sum_k n_k}{\sum_j n_j} = 1$$

Equivalently, when all $\llbracket E_k \rrbracket_s = 0$, the summation is zero. And when any $0 < \llbracket E_k \rrbracket_s < 1$, the resulting summation is a fraction between 0 and 1. Hence it is a probability.

- (c) *Concatenation:* (i) From Eq. (6):

$$\sum_i \llbracket E_1 : E_2 \rrbracket_{s_i}$$

Using Eq. (3), because s_i ranges over all Σ^* , this becomes:

$$\sum_{i,j} \llbracket E_1 \rrbracket_{s_i} \cdot \llbracket E_2 \rrbracket_{s_j}$$

This translates to:

$$\left(\sum_i \llbracket E_1 \rrbracket_{s_i} \right) \cdot \left(\sum_j \llbracket E_2 \rrbracket_{s_j} \right)$$

By the induction hypothesis, this simplifies to:

$$1 \cdot 1 = 1.$$

- (ii) Given a concatenation,

$$\llbracket E_1 : E_2 \rrbracket_s$$

By the induction hypothesis, each of E_1 and E_2 return probabilities $0 \leq p_i \leq 1$ ($i = 1, 2$). Hence their product $p_1 \cdot p_2$ must likewise be a probability.

- (d) *Kleene closure:*

- (i) Starting with Eq. (6):

$$\sum_i \llbracket E^{*f} \rrbracket_{s_i}$$

Using Eqs. (4), it translates as follows:

$$\begin{aligned} &= \llbracket E^{*f} \rrbracket_\epsilon + \sum_i \llbracket E^{*f} \rrbracket_{s_i} \quad (s_i \neq \epsilon) \\ &= (1 - f) + \sum_i f \cdot \llbracket E \rrbracket_{s'_i} \cdot \llbracket E^{*f} \rrbracket_{s_i} \quad (s_i \neq \epsilon) \end{aligned}$$

By the inductive hypothesis:

$$\sum_i \llbracket E^{*f} \rrbracket_{s_i} = (1 - f) + f \cdot 1 \cdot \sum_i \llbracket E^{*f} \rrbracket_{s_i}$$

Doing some algebraic manipulation:

$$\begin{aligned} f \sum_i \llbracket E^{*f} \rrbracket_{s_i} - \sum_i \llbracket E^{*f} \rrbracket_{s_i} &= f - 1 \\ \sum_i \llbracket E^{*f} \rrbracket_{s_i} (f - 1) &= f - 1 \\ \sum_i \llbracket E^{*f} \rrbracket_{s_i} &= 1 \end{aligned}$$

Note that the division by $f - 1$ is permitted because $f < 1$ by definition.

(ii) By induction on the length of a string s , it can be shown that

$$0 \leq \llbracket E^{*f} \rrbracket_s \leq 1$$

The base case is when $s = \epsilon$, in which case the probability is $p = 1 - f$ from the first equation in (4), and $0 < p \leq 1$. For an arbitrary $s \neq \epsilon$, the probability from the second equation in (4) is:

$$\begin{aligned} &\sum_{i=1}^n (f \cdot \llbracket E \rrbracket_{\alpha_1 \dots \alpha_{i-1}} \cdot \llbracket E^{*f} \rrbracket_{\alpha_i \dots \alpha_n}) \\ &+ f \cdot \llbracket E \rrbracket_s \cdot \llbracket E^{*f} \rrbracket_\epsilon \end{aligned}$$

By incorporating the second term into the first term's summation, this is rewritten:

$$\sum_{i=1}^n (f \cdot \llbracket E \rrbracket_{\alpha_1.. \alpha_i} \cdot \llbracket E^{*f} \rrbracket_{\alpha_{i+1}.. \alpha_n})$$

where $\alpha_m \alpha_n = \epsilon$ when $m > n$. By the inductive hypothesis over s , $\sum_i \llbracket E^{*f} \rrbracket_{\alpha_{i+1}.. \alpha_n}$ is a probability. Furthermore, by the structural induction of expressions, $\sum_i \llbracket E \rrbracket_{\alpha_1.. \alpha_i}$ is a probability. Hence their product with f is a probability.

(e) *+Closure*: Similar to (c) and (d) above. \square

3.2. Implementation of an SRE Processor

Given a regular expression, determining whether particular strings are members of its language is a tractable problem [9, 34]. There are different ways in which this may be performed. One technique is to convert the regular expression into an equivalent nondeterministic finite automaton, which can be done in polynomial time. Once this is done, a graph-searching algorithm reads a string character by character, marking states of the FA that are still eligible as paths towards an acceptance state. An advantage of the FA approach is that the nondeterministic FA can be polynomially-time translated into a deterministic FA, which will then have more efficient recognition characteristics during language recognition.

Alternatively, regular expressions can be symbolically interpreted directly. The behaviour of each regular expression operator has a corresponding operational semantics, which can be used to define a regular expression interpreter. This may be done from the perspective of either language generation or language acceptance. One technical requirement of the expression interpretation approach is that the interpreter must be able to handle the nondeterministic nature of expression derivations, since regular expressions are naturally nondeterministic in nature. The expression interpretation is similar to the FA approach, in that there is a mapping between the states of a translated FA and the derivation paths used by the interpreter when processing an expression.

Stochastic regular language recognition uses the same basic recognition schemes as conventional regular languages, with the additional requirement that probabilities be computed for strings. For example, if a FA is derived for a stochastic language, then the links are marked with probabilities. The overall probability

of accepting a given string is then computed by computing the product of all the transition probabilities used from the start state to the final accepting state. This probabilistic FA is known as a Hidden Markov Model or HMM [14]. Therefore, given a stochastic regular language as defined by SRE, the formulae of Section 3.1 are incorporated into a translated FA or a SRE interpreter.

The SRE recognition system uses the expression interpretation approach described above. The operational semantics use two relations. One relation, \longrightarrow over $E \times (\Sigma, p) \times E$, where p is a probability, represents single action transitions of expressions. This relation is denoted,

$$E \xrightarrow{(\alpha, p)} E \quad (\alpha \in \Sigma)$$

The other relation, \Longrightarrow over $E \times (\Sigma^*, p) \times E$, is the transitive closure of \longrightarrow^* , and models the generation of strings:

$$E \xrightarrow{(\alpha_1, p_1)} \dots \xrightarrow{(\alpha_k, p_k)} E \equiv E \xrightarrow{(s, p_k)} E \quad (s = \alpha_1, \dots, \alpha_k)$$

Figure 1 contains transitional rules for the relations, which define the structural operational semantics of the SRE operators [35]. These inference rules define an abstract interpreter for SRE expressions, and are the basis of an SRE recognizer. In fact, with languages such as Prolog, these rules can be compiled into Prolog statements, and then directly interpreted using Prolog's inference engine [36]. Furthermore, multiple solutions are obtained for nondeterministic SRE expressions using Prolog's builtin backtracking mechanism.

The actual implementation of the SRE processor uses the above fundamental ideas. The operational semantics implemented are a superset of the rules in Fig. 1. The implementation uses a logical grammar definition of SRE, which is part of the DCTG-GP system [37] (see Section 4). Prolog's backtracking is advantageously used to investigate different paths of an expression's derivation. In addition, string recognition is performed by pattern matching on an argument string and the generated string as shown in the transitional semantics: when a match occurs, the current derivation path is correct, while mismatches cause the current derivation to backtrack and test another possible nondeterministic path. For example, one instance of backtracking may try different terms in a Choice expression, while another may unwrap an iterative expression a varying number of times. Such backtracking is assured

$$\begin{array}{l}
\text{Action} \quad \frac{}{\alpha \xrightarrow{(\alpha,1)} \epsilon} \quad (\alpha \neq \epsilon) \\
\\
\text{Choice} \quad \frac{E_j \xrightarrow{(\alpha,p)} E'}{\sum E_i(n_i) \xrightarrow{(\alpha,q)} E'} \quad (q = p \cdot \frac{n_j}{\sum_i n_i}) \\
\\
\text{Concat}_1 \quad \frac{E \xrightarrow{(\alpha,p)} E'}{E : F \xrightarrow{(\alpha,p)} E' : F} \\
\\
\text{Concat}_2 \quad \frac{F \xrightarrow{(\alpha,p)} F'}{\epsilon : F \xrightarrow{(\alpha,p)} F'} \\
\\
\text{Kleene}_1 \quad \frac{}{E^*f \xrightarrow{(\epsilon,1-f)} \epsilon} \\
\\
\text{Kleene}_2 \quad \frac{E \xrightarrow{(\alpha,p)} E'}{E^*f \xrightarrow{(\alpha,p \cdot f)} E' : E^*f} \quad (\alpha \neq \epsilon) \\
\\
+\text{Closure} \quad \frac{E \xrightarrow{(\alpha,p)} E'}{E^+f \xrightarrow{(\alpha,p)} E' : E^*f} \quad (\alpha \neq \epsilon) \\
\\
\text{Strings}_1 \quad \frac{E \xrightarrow{(\alpha,p)} \epsilon}{E \xrightarrow{(\alpha,p)} \epsilon} \\
\\
\text{Strings}_2 \quad \frac{E \xrightarrow{(\alpha,p)} E' \quad E' \xrightarrow{(s,q)} E''}{E \xrightarrow{(\alpha s, pq)} E''}
\end{array}$$

Figure 1. Transitional semantics of SRE.

of terminating because of the finite size of input strings to be checked, as well as the assertion within the SRE semantics that empty strings ϵ can never be generated within the generative component of iterative operators (they can only be generated when the iteration terminates).

One advantage of a stochastic language is that the computed probabilities of strings can be used as an efficiency mechanism during expression recognition. The implementation of the SRE recognizer is such that the probability of intermediate strings are always known

throughout the interpreter. When the current probability becomes smaller than a user-supplied threshold, the current derivation path can be forced to terminate. This prunes derivations of an expression which yield probabilities too small to be of consequence. Of course, setting this threshold too large results in inaccurate probability values for recognized strings, and may even erroneously reject legal strings. However, for many experiments, especially with large strings to be recognized, this speeds up processing significantly.

4. Genetic Programming System

4.1. Grammatical SRE and gSRE

The GP system used for the SRE experiments is the DCTG-GP system [37]. DCTG-GP performs grammar-based genetic programming, in which the target language of the evolved program population is defined in terms of a context-free grammar [26, 38–40]. A major advantage of grammatical GP systems is that the search space is syntactically constrained so that evolution is given a helpful push towards program structures that are more sensible for the problem at hand.

The grammar used by DCTG-GP is a definite clause translation grammar, or DCTG [41]. A DCTG is a logical version of a context-free attribute grammar. Each DCTG production has a syntactic component, which defines a context-free syntax rule. In addition, each production can have included with it one or more semantic components. A semantic component defines some characteristic of the syntactic component to which it is attached. For example, one important SRE characteristic that is defined in the DCTG grammar is the string recognition algorithm of Section 3.2. During fitness evaluation, the GP system tests whether gSRE expressions can recognize different example strings of the target language. Hence the operational semantics of gSRE are encoded so that expression interpretation attempts to recognize the membership of strings, and produce their corresponding probabilities if so recognized. Given a string to recognize, the actual implementation finds the probability of the largest prefix recognized by an expression (more details are in Section 5). Since the operational semantics of the SRE operators are modular in nature, their recognition behaviours can be encoded with the grammar rules that define the syntax of the operators themselves. The overall result of this is a compact definition of the SRE language, in which the syntax and semantics are conveniently unified together.


```

choice ::= guardedexpr^A1, intval^B1, guardedexpr^A2, intval^B2
<:>
(recognize(S, S2, Sum, PrSoFar, Pr) :-
  B1^construct(Val),
  Pr2 is (Val/Sum)*PrSoFar,
  A1^recognize(S, S2, Pr2, Pr)),
(recognize(S, S2, Sum, PrSoFar, Pr) :-
  B2^construct(Val),
  Pr2 is (Val/Sum)*PrSoFar,
  A2^recognize(S, S2, Pr2, Pr)).

```

Figure 2. Operational semantics of choice operator (excerpt).

To give a flavour of the system, Fig. 2 shows an excerpt of the DCTG grammar and semantics for the guarded choice operator. The rules for `guardedexpr` (not shown) defining the guarded choice terms are encoded within the grammar as syntactic constraints. Rather than permitting any SRE expression as a term within a choice operator, only uniquely guarded terms are permitted. This generally forces expressions to be more concise, with no loss in descriptiveness. The rule in Fig. 2 is pertinent when there are two choices possible. The rule `recognize(S, S2, Sum, PrSoFar, Pr)` has 5 arguments: the string S at the start of processing of this choice expression; the string $S2$ after processing ($S2$ is either equal to S or a suffix of it); the Sum of the probability values terms in the choice list (i.e. the overall denominator value); the probability $PrSoFar$ computed so far while processing the current string (other expression components may have read earlier prefixes of the example string, and have this probability); and the final computed probability Pr after processing S is completed.

There are two rules for `recognize`, and each rule processes one of the terms from the pair of choice expressions. The call to `construct` retrieves the actual integer value from the probability field for that term. The probability for that term is calculated, and multiplied by the overall probability so far. This new intermediate probability is then passed to the recognition semantics for the expression embedded in that choice term. During processing, both rules will be invoked—the first followed by the second—because there are multiple ways a string can be recognized by an sGRE expression. The operational semantics for gSRE will exhaustively try all rules until the string is completely recognized, and no alternative derivations of the gSRE expression are possible. All the probabilities obtained

for these difference derivation paths are collected and summed, to yield an overall probability for that string (or the greatest prefix read, as the case may be).

Another syntactic constraint applied to both SRE and gSRE in the experiments in Section 5 is the following. Although not specified in the grammar of SRE (or gSRE), the grammatical definition of SRE disallows iterative operators to be directly nested within one another. In other words, expressions such as

$$((E)^*f)^*f \quad \text{or} \quad (((E)^*f)^+f)^*f)^+f$$

are not allowed. The reason for this restriction is a pragmatic one. When GP was performed without this restriction, many programs had multiply nested iterative expressions. Such expressions are relatively expensive to interpret, due to the variety of nondeterministic paths possible for interpreting them. In addition, nested iteration typically results in strings with very low probabilities, since there is a probabilistic factor f associated with executing every nested iterative expression. Moreover, the expense of nested iteration is not justified by results, since any of these expressions can be replaced with a semantically equivalent expression that uses only one iterative operator. This restriction does *not* imply that an expression like

$$((a : E^*f)^*f$$

is illegal, since the concatenation operator means that the iteration operators are not directly nested.

4.2. Other GP System Details

DCTG-GP uses standard GP strategies, such as tournament or roulette-wheel selection, and steady-state or

generational evolution. Relevant experimental parameters will be illustrated in Section 5. The system is implemented in Sicstus Prolog 3 on both Windows 98 and Silicon Graphics platforms.

5. Experiments

5.1. General Strategy

The inference of a stochastic language can be considered to involve two different objectives. Given a training set of positive (and possibly negative) examples, one task is to infer a language which correctly classifies the training examples. This is equivalent to non-stochastic language inference. An additional task required for stochastic language inference, however, is to ascertain the stochastic distribution of the training examples. One might naively presume that a statistical analysis of the training set could be performed, and the results applied to the inferred language. Unfortunately, the situation is typically more complicated than this, because the representation of the stochastic language as used in the hypothesis will not likely permit a straight-forward application of the final string distributions to its internal encoding. For example, if an HMM representation is used in hypotheses, finding appropriate probability values for intermediate links in the network that will correspond to the example set distribution is a challenging task. The significance of the problem of determining distributions for HMM's and context-free languages has spawned specialized training algorithms [7, 8].

The inference strategy undertaken with the GP experiments is to let evolution determine stochastic distributions in concert with example classification. Since SRE incorporates probability values directly in expressions, treating numeric probability fields is straight-forward in GP. It was found that this approach was sufficient for many experiments undertaken. In fact, it was discovered that evolution using local search for fine-tuning probability parameters lent no advantage over simple evolution of the parameters.

The training sets used in the GP experiments consist of sets of positive examples for the target language to be inferred. Each member of the set is a string, along with its frequency with respect to the total number of strings in the set (typically 1000). Since the format of the target languages is already known via a stochastic regular expression or grammar, generating these sets is

straight-forward. Unlike conventional language inference, the implicit probability distributions in training example sets permits stochastic languages to forgo the need for negative examples. This is because the inference of a distribution that matches that of the training set will automatically account for 'negative examples', which have 0 probability in the distribution.

Stochastic language inference incorporates an implicit degree of error in any inferred solution. This has ramifications on the GP fitness evaluation described below. It also can be used to boost efficiency of computations performed during inference. As detailed in Section 3.2, string recognition can be pre-empted when intermediate probabilities become smaller than some threshold limit set for the experiment. Similarly, the test set can be pruned of strings whose frequency is below some limit set by the user. This limit parameter should be set with the recognition threshold in mind. For example, if the threshold is set to 0.001, then the test set limit could be likewise set to 1 for a test set of size 1000. Of course, there may be many nondeterministic derivations of an expression when recognizing a string, and all the probabilities of these derivations will be summed to an overall probability for that string. The less discriminating the recognition threshold and test set limit, the more precise (albeit slower) the results.

Since GP experiments use a steady-state algorithm, there are not any discrete generations. For convenience, however, a new generation is said to have occurred every K reproductions, where K is the population size. Between generations, the test set is regenerated. This prevents overfitting to one set of test data, and reflects the nature of the stochastic languages, as each test set reflects a sampling from the actual distribution. One disadvantage, however, is that a discrete test set is an approximation of the real distribution of the language, and hence this introduces an unavoidable measure of noise. This noise is compensated by the fact that multiple test sets are used during successive generations, and their cumulative effect should reflect a more accurate model of the target distribution. However, the population is not reevaluated for each newly generated test set, and so the fitnesses of much of the population may be legacy values from earlier generations. This is acceptable, because the test sets used for those generations are presumed to be as statistically valid as those from any other generation.

The fitness evaluation strategy used in the experiments is a modified χ^2 test [42]. The known distribution is taken to be the set T of test examples, and the

experimental set will be the results of the SRE recognition algorithm on each member $t_i \in T$. Each test set example string is given to the SRE processor, and an overall probability $Pr(t_i)$ for that string is computed. Non-membership is reflected in a probability of 0. The fitness formula is:

$$\text{Fitness} = \sum_{t_i \in T} \begin{cases} \frac{(d_i - (Pr(t_i) \times N))^2}{d_i} & Pr(t_i) \geq 0 \\ \frac{|t_i| - |maxpref_i|}{|t_i|} \cdot d_i & Pr(t_i) = 0 \end{cases}$$

where d_i is the frequency of example t_i in test set T , $N = |T|$, and $maxpref_i$ is the maximum prefix of t_i recognized. The first term is the χ^2 formula, and it is used when the example string t_i is completely recognized. The second formula is used when only a prefix of t_i is recognized. Its value is inversely proportional to the size of the prefix recognized. Should none of t_i be recognized, then this value becomes $2 \cdot d_i$ (a normal χ^2 formula would use just d_i). This prefix scoring gives credit to expressions that recognize portions of the examples, which helps drive evolution towards expressions that recognize complete examples.

5.2. Experiment 1: Stochastic Iteration

The first experiment uses a simple stochastic regular language which can be naturally encoded in SRE. The main intention of this experiment is to test the evolvability of stochastic Kleene closure as modeled in SRE. The target language is a stochastic rendition of a regular language suggested by Tomita¹ from his popular benchmarks for machine learning [20]. The target language written in SRE is:

$$L_1 = a^{*.5} : b^{*.5} : a^{*.5} : b^{*.5}$$

This is a non-trivial language, especially in the stochastic domain, as the overall distribution of each a and b term in all the strings should conform to the given probability of 0.5. These terms may also generate empty strings, should iterations terminate immediately.

The parameters for the experiment are in Table 1. Most are self-explanatory, and the fitness function strategy was discussed earlier in Section 5.1. The initial population is oversampled, and the running population is pruned from it using tournament selection. Replacement is done using a reverse tournament selection (a sample of K members are randomly selected,

Table 1. Parameters (L_1).

Parameter	Value
Target language	gSRE
Terminals	a, b
Fitness function	Modified χ^2
Generation type	Steady-state
Initial population size	750
Running population size	500
Unique population members	Yes
Maximum generations	50
Probability of crossover	0.90
Probability of mutation	0.10
Probability internal crossover	0.90
Probability terminal mutation	0.75
Probability numeric mutation	0.50
Numeric mutation range	± 0.1
Max reproduction attempts	3
Initial population shape	Ramped half & half
Min/max depth initial popn.	6, 12
Max depth offspring	24
Tournament size	5
Test set size	1000
Max test string size	Approx. 20
Min test example frequency	3
SRE probability limit	0.0001

and the member with the lowest fitness is selected to be replaced).

Mutation is performed on either terminal or nonterminals. If a nonterminal is to be mutated, there is a 0.5 probability that it should be a numeric field. When a numeric field is selected for mutation, its current value is perturbed $\pm 10\%$ of the entire range for that numeric type (a range of ± 100 for integers, and ± 0.1 for probabilities).

A test set is generated before every generation. Initially, 1000 strings are generated for L_1 , and their frequencies are tallied. The maximum string size is approximately 20 (some may exceed this length). Should there be less than 3 instances of a given string, it is pruned from the test set. This means that there are typically between 55 to 60 unique strings in the test set, each of which has its particular frequency for that particular sample of the language. The number of unique strings in the test set is important for χ^2 analyses, as it is equivalent to bin size in the χ^2 formula.

Table 2. Summary L_1 .

Total runs	15
# unique examples	65
Avg. test set χ^2 (50 cases)	142.22
Fitness	
Min	89.4 ($\chi^2 = 88.8$)
Max	251.55 ($\chi^2 = 203.75$)
Avg	127.91 ($\chi^2 = 124.27$)

Summary statistics for the best solutions from 15 runs are given in Fig. 2. These values are obtained using a common test set, since each run will have used a different test set during its prior evolution. An average χ^2 test of the test sets themselves is included, in order to better evaluate the expression results. 50 pairs of random test sets were generated. One of the pair was fixed as the independent variable, while the other was the dependent variable. The sets were filtered for frequencies below the minimal test example frequency (2 in Fig. 1), and the χ^2 was computed. The resulting 50 χ^2 values were averaged.

A performance chart of the best and average population fitness averaged for 15 runs is in Fig. 3. It can be

seen that convergence to a local optimum has largely occurred by generation 10.

The best solution found ($\chi^2 = 88.8$), is:

$$a^{*.5} : b^{*.44} : a^{.49} : b^{.49}$$

This is a nearly perfect solution, and the iterative probabilities within the range of what might be expected given the stochastic error inherent with the random test sets. The second best solution ($\chi^2 = 89.63$) is:

$$a^{*.52} : b^{*.49} : a^{.37} : (a(94) + b(759))^{*.54}$$

The last term is interesting, in that the erroneous choice of a is not too acute a problem, given the low probability of choosing it (0.11).

One of the poorer solutions ($\chi^2 = 132.85$) is:

$$(a(468) + b(235))^{*.25} : a^{*.55} : (a(182) + b(963))^{*.73}$$

The inaccuracy occurs with the first term, which erroneously permits b to occur too frequently, even though the low probability of 0.25 for the enclosing iteration helps reduce its likelihood.

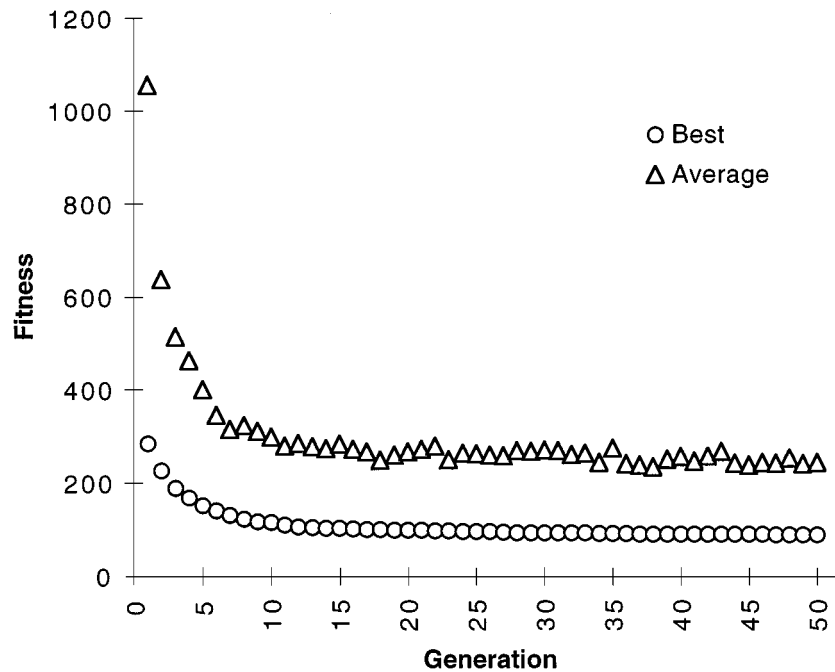


Figure 3. Fitness curves (avg. 15 runs).

Table 3. Summary L_2 .

Total runs	50
# unique examples	35
Avg. test set χ^2 (50 cases)	99.75
Fitness	
Min	66.39 ($\chi^2 = 65.06$)
Max	281.93 ($\chi^2 = 274.1$)
Avg	199.5 ($\chi^2 = 193.26$)

- $S \rightarrow a S \quad (0.2)$
- $S \rightarrow b A \quad (0.8)$
- $A \rightarrow a B \quad (0.7)$
- $A \rightarrow b S \quad (0.3)$
- $B \rightarrow a A \quad (0.4)$
- $B \rightarrow b B \quad (0.1)$
- $B \rightarrow \epsilon \quad (0.5)$

Figure 4. Target language.

The worst solution obtained ($\chi^2 = 203.75$) is:

$$\begin{aligned}
 &(a^{*.58} : (a(362) + b(805))^{+.67})^{+.04} \\
 &: (b^{+.2} : b^{+.59} : ((a : (a(364) + b(320))^{*.08}(320) \\
 &+ b(947) : (a(191) + b(141))^{+.67})^{+.2} \\
 &: ((a(352) + b(360))^{+.56} : b^{+.59})^{*.04})^{*.04}
 \end{aligned}$$

Note the repetition of particular numeric fields, such as 320 and 0.04, which is a sign of population convergence. Simplifying this expression by removing iterative probabilities less than 0.10 and expanding +Closure terms, it becomes:

$$a^{*.58} : (a(362) + b(805))^{+.67}$$

which is obviously a suboptimal solution. This example shows the nature of introns within SRE expressions: virtually any expression can be intron code, so long as the associated choice or iterative probability is low enough.

5.3. Experiment 2: Stochastic Regular Grammar

The second experiment evolves a more complex stochastic regular language. The target language L_2 is taken from [10], and is defined by the stochastic regular grammar in Fig. 4. Each production has a probability on the right, which denotes the probability that rule is selected with respect to the other productions for that nonterminal.

The experimental parameters for these runs are identical to those in Fig. 1. The summary for 50 runs are in Fig. 3. A performance plot for the best fitness and average population fitness averaged for the 50 runs is in Fig. 5.

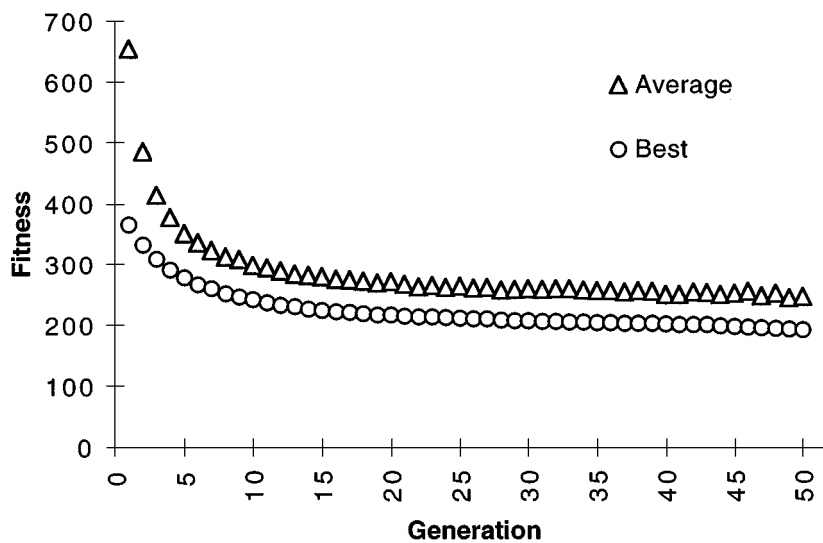


Figure 5. Fitness curves (avg. 50 runs).

The best solution ($\chi^2 = 65.06$) is:

$$\begin{aligned} & ((a : (a(674) + b(895))^{*.03})(895) + b : b(960))^{*.37} \\ & : (b : (a : (a)^{+.02}(674) + b : b(284))^{*.29}) \\ & : (a : a^{+.02}(895) + b : b(284))^{*.29} \\ & : b : (a(490) + b(524))^{*.03} : b^{*.09} \\ & : b : a^{+.02} : (a : a^{+.02}(895) + b(284))^{*.37} \end{aligned}$$

Simplifying by expanding + Closures and removing terms with probabilities less than 0.03, this becomes:

$$\begin{aligned} & (a : (a(674) + b(895))(895) + b : b(960))^{*.37} \\ & : (b : (a : a(674) + b : b(284))^{*.29}) \\ & : (a : a(895) + b : b(284))^{*.29} : b^{*.09} \\ & : b : a : (a : a(895) + b(284))^{*.37} \end{aligned}$$

It is difficult to see how this expression maps to the target grammar of Fig. 4, and an intuitive mapping may not even exist. However, its χ^2 is impressive compared to the test set average.

5.4. Limitations

Many language inference algorithms are easily thwarted by target languages having characteristics antagonistic to the peculiarities of the algorithm in question. Often, these languages are only subtly different from ones that the algorithms have no problems inferring.

The GP paradigm suffers a similar limitation. A variation of the language in Section 5.3 was tried,

$$L'_2 = L_2(9) + bbaaabab(1)$$

which is just language L_2 with an additional string $bbaaabab$ with probability 10%. 50 runs were performed using the same parameters as Fig. 1. None of the runs found an acceptably close solution: the best solution had a fitness of 259 and $\chi^2 = 234$ (bin size = 40).

One reason that GP had problems evolving L'_2 can be attributed to the linguistic characteristics of SRE. Even though the above definition of L'_2 is a concise statement of the language, the evolutionary process tries to unify the term $bbaaabab$ and L_2 together in a regular expression. This is difficult to do, because this string is an anomaly with respect to the other strings in L_2 . Considering the stochastic regular grammar used to generate

L_2 , it is clear that strings are derived progressively and incrementally from one another, and so strings of L_2 equal in length to $bbaaabab$ are natural extensions of smaller strings of the language. The anomalous string, however, is not derivable from L_2 , and hence a natural model of the union of these languages in SRE cannot be inferred. This is especially true given that $bbaaabab$ has a 10% probability, which makes it a populous member. If it had a smaller probability, it might be ignored as noise.

The above must be considered in light of the linguistic nature of all formal languages: some representations more naturally model particular languages than others. Even though regular expressions, finite automata and regular grammars have the same expressive power, some languages are more naturally and concisely denoted by regular expressions than by finite automata, and vice versa. It could be the case that another representation language, for example HMM's, may more naturally denote L'_2 than SRE.

6. Conclusion

This paper presented a new means for evolving stochastic regular languages. Using a probabilistic version of regular expressions as a language for evolution, genetic programming is capable of evolving accurate expressions for stochastic regular languages. However, some stochastic regular languages are more amenable to successful evolution than others. It can be speculated that languages in which members have structural similarities with one another are the most suitable for this paradigm. For more complex languages, more sophisticated evolutionary techniques may be required.

It was found during experimentation that SRE had no evolutionary advantage over gSRE with respect to the quality of solutions discovered. On the other hand, SRE expressions were less efficient to process, and runs took much longer than the gSRE ones.

The use of SRE in a genetic programming context presents advantages over other evolutionary experiments with stochastic languages. One advantage is that SRE is akin to a programming language, with operators that have syntactic and semantic definitions akin to conventional languages. Since GP is typically applied towards such languages as Lisp, the encoding and processing of SRE within a GP environment is straight-forward. More importantly, however, is that SRE has linguistic advantages over finite automata and regular grammars: some stochastic languages are more

naturally encoded in SRE than these other representations. The L_1 experiment is a clear example of this point. The linguistic clarity of L_2 is less apparent, although the solution is not overly complex compared to the target grammar.

Like [23], this work uses a regular expression language directly for GP. His work required fairly large populations and parallel populations in order to evolve the Tomita languages. The fitness strategy used here is similar to that used in [25, 30], in that both language recognition performance and prefix consumption are taken into consideration.

There are many directions for future work. The GP strategies used here were fairly conventional, and more sophisticated approaches may be more applicable to stochastic languages. In the experiments, the wide degree of qualitative variations between runs indicates that evolution quickly gets stuck at suboptimal solutions. Parallel subpopulations may help in this regard. Although it was found that local search using hill-climbing over numeric fields was not advantageous to evolution, it is worth investigating the utility of more sophisticated local search techniques akin to those used in stochastic context-free languages (eg. the inside-outside algorithm).

Currently, the applicability of SRE in bioinformatics problems is being investigated. A fundamental problem in DNA and protein sequencing is to determine a common pattern shared amongst a family of sequences [43], which can be used for both search and analytical purposes. A number of techniques, such as HMM's and regular pattern languages, are used for this purpose. SRE is a natural vehicle for this problem area, since its regular expression basis conforms to the pattern languages commonly used (eg. that used in the PROSITE database [44]), while its stochastic features conveniently model the probabilistic characteristics of DNA sequences themselves.

Acknowledgment

Thanks to Tom Jenkyns for helpful discussions about probability. This research is supported by NSERC Operating Grant 138467-1998.

Note

1. His language is $a^*b^*a^*b^*$.

References

1. D. Angluin, "Computational learning theory: survey and selected bibliography," in *Proceedings of the 24th Annual ACM Symposium on the Theory of Computing*, ACM Press: New York, 1992, pp. 351–369.
2. K.S. Fu and T.L. Booth, "Grammatical inference: introduction and survey—Part I," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 5, no. 1, pp. 95–111, January 1975.
3. K.S. Fu and T.L. Booth, "Grammatical inference: introduction and survey—Part II," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 5, no. 4, pp. 409–423, July 1975.
4. Y. Sakakibara, "Recent advances of grammatical inference," *Theoretical Computer Science*, vol. 185, pp. 15–45, 1997.
5. K.S. Fu and T. Huang, "Stochastic grammars and languages," *International Journal of Computer and Information Sciences*, vol. 1, no. 2, pp. 135–170, 1972.
6. K.S. Fu, *Syntactic Pattern Recognition and Applications*, Prentice-Hall: Englewood Cliffs, NJ, 1982.
7. K. Lari and S.J. Young, "The estimation of stochastic context-free grammars using the Inside-Outside algorithm," *Computer Speech and Language*, vol. 4, pp. 35–56, 1990.
8. E. Charniak, *Statistical Language Learning*, MIT Press: Cambridge, MA, 1993.
9. J.E. Hopcroft and J.D. Ullman, *Introduction to Automata Theory, Languages, and Computation*, Addison Wesley: Reading, MA, 1979.
10. R.C. Carrasco and M.L. Forcada, "Inferring stochastic regular grammars with recurrent neural networks," in *Proceedings 3rd International Colloquium on Grammatical Inference (ICGI96)*, edited by I. Miclet and C. de la Huguera, Springer-Verlag: Berlin 1996, pp. 274–286. LNAI 1147.
11. R.C. Carrasco and J. Oncina, "Learning deterministic regular grammars from stochastic samples in polynomial time," Technical Report DLSI-96-01, Universidad de Alicante, April 1998.
12. F.J. Maryanski and T.L. Booth, "Inference of finite-state probabilistic grammars," *IEEE Transactions on Computers*, vol. C26, pp. 521–536, 1977.
13. A. van der Mude and A. Walker, "On the inference of stochastic regular grammars," *Information and Control*, vol. 38, pp. 310–329, 1978.
14. L.R. Rabiner and B.H. Juang, "An introduction to Hidden markov models," in *IEEE ASSP Magazine*, pp. 4–16, January 1986.
15. L.R. Rabiner, "A tutorial on Hidden markov models and selected applications in speech recognition," in *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, February 1989.
16. V.K. Garg, R. Kumar, and S.I. Marcus, "Probabilistic language framework for stochastic discrete event systems," Technical Report 96-18, Institute for Systems Research, April 1996. <http://www.isr.umc.edu/>.
17. H. Zhou and J.J. Grefenstette, "Induction of finite automata by genetic algorithms," in *Proc. 1986 IEEE Intl. Conference on Systems, Man, and Cybernetics*, Atlanta, GA, IEEE Press: New York, 1986, pp. 170–174.
18. B.D. Dunay, F.E. Petry, and B.P. Buckles, "Regular language induction with genetic programming," in *Proc. 1st IEEE Conference on Evolutionary Computation*, 1994, pp. 396–400.
19. P. Dupont, "Regular grammatical inference from positive and negative samples by genetic search: the GIG method," in *2nd*

- Intl. Coll. on Grammatical Inference and Applications*, Springer-Verlag: Berlin, 1994, pp. 236–245.
20. M. Tomita, “Dynamic construction of finite automata from examples using hill-climbing,” in *4th Annual Conf. of the Cognitive Science Soc.*, 1982, pp. 105–108.
 21. S. Brave, “Evolving deterministic finite automata using cellular encoding,” in *Proc. Genetic Programming 1997*, Stanford University, CA, USA, edited by J.R. Koza et al., Morgan Kaufmann: Los Altos, CA, 1997, pp. 39–44.
 22. T. Longshaw, “Evolutionary learning of large grammars,” in *Proc. Genetic Programming 1997*, Stanford University, CA, USA, edited by J.R. Koza et al., Morgan Kaufmann, Los Altos, CA, 1997, pp. 406–409.
 23. B. Svingen, “Learning regular languages using genetic programming,” in *Proc. Genetic Programming 1998*, edited by J.R. Koza et al., Morgan Kaufmann: Los Altos, CA, 1998, pp. 374–376.
 24. P. Wyard, “Context free grammar induction using genetic algorithms,” in *Proceedings 4th International Conference on Genetic Algorithms*, 1991, pp. 514–518.
 25. M.M. Lankhorst, “Grammatical inference with a genetic algorithm,” in *Proceedings of the 1994 EUROSIM Conference on Massively Parallel Processing Applications and Development*, 1994, pp. 423–430.
 26. S. Lucas, “Structuring chromosomes for context-free grammar evolution,” in *Proceedings 1st International Conference on Evolutionary Computation*, IEEE Press: New York, 1994, pp. 130–135.
 27. S. Sen and J. Janakiraman, “Learning to construct pushdown automata for accepting deterministic context-free languages,” *Applications of Artificial Intelligence X: Knowledge-Based Systems*, vol. 1707, pp. 207–213, 1992.
 28. M.M. Lankhorst, “A genetic algorithm for the induction of push-down automata,” in *1995 IEEE Intl. Conference on Evolutionary Computation*, IEEE Press: New York, 1995.
 29. B.D. Dunay and F.E. Petry, “Solving complex problems with genetic algorithms,” in *Proc. 6th Intl. Conf. on Genetic Algorithms*, edited by L. Eshelman, Morgan Kaufmann: Los Altos, CA, 1995.
 30. M. Schwehm and A. Ost, “Inference of stochastic regular grammars by massively parallel genetic algorithms,” in *Proc. 6th Intl. Conf. on Genetic Algorithms*, Morgan-Kaufmann: Los Altos, CA, 1995.
 31. T.E. Kammeyer and R.K. Belew, “Stochastic context-free grammar induction with a genetic algorithm using local search,” in *Foundations of Genetic Algorithms IV*, edited by R.K. Belew and M. Vode, Morgan-Kaufmann: Los Altos, CA, 1997.
 32. J. Stoy, *Denotational Semantics*, MIT Press: Cambridge, MA, 1977.
 33. K. Subrahmaniam, *A Primer in Probability*, Marcel Dekker: New York, 1979.
 34. M. Sipser, *Introduction to the Theory of Computation*, PWS Pub. Co., 1996.
 35. M. Hennessy, *The Semantics of Programming Languages—An Elementary Introduction Using Structural Operational Semantics*, John Wiley and Sons: New York, 1990.
 36. W.F. Clocksin and C.S. Mellish, *Programming in Prolog*, 4th ed, Springer-Verlag: Berlin, 1994.
 37. B.J. Ross, “Logic-based genetic programming with definite clause translation grammars,” in *Proc. GECCO-99*, edited by J.R. Koza et al., 1999.
 38. P.A. Whigham, “Grammatically-based genetic programming,” in *Proceedings Workshop on Genetic Programming: From Theory to Real-World Applications*, edited by J.P. Rosca, 1995, pp. 31–41.
 39. M.L. Wong and K.S. Leung, “Learning programs in different paradigms using genetic programming,” in *Proceedings 4th Congress of the Italian Association for AI*, 1995, pp. 353–364.
 40. A. Geyer-Shulz, “The next 700 programming languages for genetic programming,” in *Proc. Genetic Programming 1997*, Stanford University, CA, USA, edited by J.R. Koza et al., Morgan Kaufmann: Los Altos, CA, 1997, pp. 128–136.
 41. H. Abramson and V. Dahl, *Logic Grammars*, Springer-Verlag, 1989.
 42. W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery, *Numerical Recipes in C*, 2nd edn. Cambridge University Press: Cambridge, 1992.
 43. A. Brazma, I. Jonassen, I. Eidhammer, and D. Gilbert, “Approaches to the automatic discovery of patterns in biosequences,” Technical Report 113, Department of Informatics, University of Bergen, Norway, December 1995.
 44. K. Hofmann, P. Bucher, L. Falquet, and A. Bairoch, “The PROSITE database, its status in 1999,” *Nucleic Acids Research*, vol. 27, no. 1, pp. 215–219, 1999.