

Probabilistic Regular Grammar Inference Algorithm Using Incremental Technique

Torsak Penpinun¹⁺ and Athasit Surarerks²

¹ pentorja@gmail.com, Department of Computer Engineering, Faculty of Engineering, Chulalongkorn University, Bangkok, Thailand

² athasit.s@chula.ac.th, Department of Computer Engineering, Faculty of Engineering, Chulalongkorn University, Bangkok, Thailand

Abstract. Grammatical inference has been studied for a long time where grammar is illustrated by a collection of re-writing rules, together with their probabilities. We are interested in regular language model which can be recognized by a finite state machine. The most popular technique is an Alergia algorithm. The objective is to construct a probabilistic finite state machine using only positive examples together with their probabilities (or frequency). In this work, we introduce a probabilistic grammatical inference algorithm in order to construct a prefix tree. The algorithm starts by considering the shortest positive example. Two types of regular grammar rules (productions) are introduced. Our experimental results show that the probabilities obtained from our probabilistic finite state machine can be more accurate than the one obtained from the previous algorithm. We hope that our algorithm will be an alternative way for constructing a probabilistic finite state machine.

Keywords: grammar inference, probabilistic finite state machine, incremental technique.

1. Introduction

Grammatical inference or grammatical induction has been studied for a long time. The theoretical foundations in grammatical inference were first introduced by M.E. Gold [1]. This topic can be considered a problem in natural language processing (NLP) and in human cognition. A grammar is described as a model of sentence-level phenomena in language. Many researches concentrated on how to learn a grammar based on Chomsky hierarchy [2] from some positive examples of words in specific language. The typically classes of grammars can be found [1, 3, 4, and 5]. Grammatical inference can also be applied in many applications such as automatic speech recognition, statistical machine translation and information retrieval, etc.

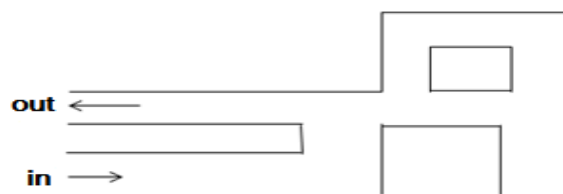
Considering the structural pattern recognition, and computational linguistic, the important problem in inference is that of dealing with positive and negative examples. In detail, grammar can be illustrated by a collection of re-writing rules, together with their probabilities. Such rules (or production rules) are used for producing words in some specific languages. The research topic is to determine production rules which can generate most of words in the target language.

Many research domains including pattern recognition, the process of learning is a set of production rules for strings in a formal language. The rules used to describe how to generate strings from alphabet that are valid according to the syntax. Normally use for describing a language model which can represent a set of sequences of processes. There are limited literatures recognized a language from examples, a probabilistic finite automaton is usually used for an inference model. Many researches interest in algorithm to inference model but Alergia algorithm [6], the famous one, is more efficient and accurate. After we studied Alergia algorithm, we have found that it is a complex grammatical inference and takes a lot of computing time.

⁺ Corresponding author.
E-mail address: pentorja@gmail.com.

In this work, we concentrated on how a language can be recognized using some positive examples. We scoped only on the regular languages. Our concept is to construct a grammatical inference algorithm with positive examples. The technique is to represent a language using grammar model together with its probability comparing with the Alergia algorithm. We then show in our work by some examples that the probability obtained from our algorithm are closed to the original comparing to the previous work.

2. Background Knowledge



Example 2: Given a regular grammar G defined by $\{S \rightarrow aS, S \rightarrow bT, T \rightarrow aT, T \rightarrow \lambda\}$. Let L be a language that every string in the language can be produced from G . It can see that a string baa is in L since “ baa ” can be produced by starting from S ,

$S \Rightarrow bT$	using the production $S \rightarrow bT$,
$S \Rightarrow baT$	using the production $T \rightarrow aT$,
$S \Rightarrow baaT$	using the production $T \rightarrow aT$,
$S \Rightarrow baa$	using the production $T \rightarrow \lambda$.

Probabilistic Finite State Machine (PFSM)

Finite state machine is a mathematical model for recognizing a regular language with respects to the input. We will recall now a probabilistic model for a finite state machine as follows:

Definition 4: A probabilistic finite state machine is composed of

- Q A finite set of states where q_0 is the initial state,
- Σ A finite alphabet is a set of all characters,
- T A transition function where $T:(Q \times \Sigma \cup \{\lambda\}) \rightarrow (Q \times p)$ where p is the probability of the transition and λ is an empty string.

The probability that a string $w = x_1x_2 \dots x_n$ (where x_i is in Σ) is in the language associated to the PFSM is the product of p 's of every transition for w .

Example 3: Let PFSM be a probabilistic finite state machine as shown in Fig 2. The probability of an empty string (λ) is $1/2$. The probability of "ab" can be computed by $1/4 \times 1/4 \times 2/3 = 1/24$.

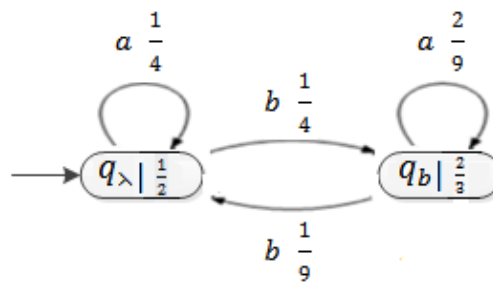


Fig. 2: An example of a probabilistic finite state machine with two states.

Alergia Algorithm and Hoeffding bound

The Alergia algorithm [6, 7], is a learning algorithm which works by merging the states of a generated finite state machine from a probabilistic perspective. The algorithm starts by generating a Prefix Tree Acceptor (PTA) from the input (positive) examples and analyses the relative frequency of outgoing transactions at every node. The algorithm compares probabilities between two states using the Hoeffding bound concept [8].

Alergia Algorithm

Let Σ be the finite set of alphabet $\{a, b\}$.

Let S be the samples set of strings

Step 1: Initialize Prefix Tree Acceptor

A is a stochastic prefix tree acceptor from S

Step 2: Compare two successor nodes (start from the root of A)

Using Hoeffding bound between successor (first node(A) to last node(A))

Step 3: Merging techniques

If Hoeffding satisfied,

Merge (firstnode(A), the last node(A))

Determinize(A)

Repeat Step 3 until no successor.

Step 4: The result

Return A

Example 4: From PFSM defined in Example 7. Given some positive examples as illustrated in Table 1, the PFSM obtained from the Alergia algorithm contains of two states as shown in Fig. 3.

Table 1: Input positive examples

Strings	Frequency	Probability	Strings	Frequency	Probability
λ	490	0.49	abaa	2	0.002
a	128	0.128	abab	2	0.002
b	170	0.17	abba	2	0.002
aa	31	0.031	abbb	1	0.001
ab	42	0.042	baaa	2	0.002
ba	38	0.038	baab	2	0.002
bb	14	0.014	baba	1	0.001
aaa	8	0.008	babb	1	0.001
aab	10	0.01	bbaa	1	0.001
aba	10	0.01	bbab	1	0.001
abb	4	0.004	bbba	1	0.001
baa	9	0.009	aaaaa	1	0.001
bab	4	0.004	aaaab	1	0.001
bba	3	0.003	aaaba	1	0.001
bbb	6	0.006	aabaa	1	0.001
aaaa	2	0.002	aabab	1	0.001
aaab	2	0.002	aabba	1	0.001
aaba	3	0.003	abbba	1	0.001
aabb	2	0.002	abbab	1	0.001

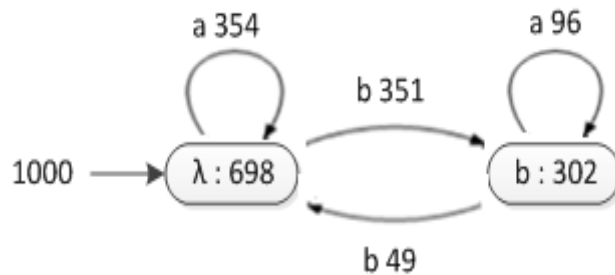


Fig. 3: Probabilistic finite state machine generated by Alergia algorithm.

3. Our Works

Alergia algorithm is always a most popular technique for learning PFSM, but the problem is the number of states of the PFSM are probably large. The study of Xuanyi Qi [9] aimed to analyse and improve the algorithm via minimization of deterministic finite automaton. This can reduce the number of states of the output PFSM. We remark that time complexity of the algorithm is still a problem. In this work, we propose an alternative way to construct a probabilistic grammar using incremental technique. The result of PFSM can be obtained from such productions.

Algorithm: Probabilistic grammar inference

Let Σ be the finite set of alphabet.

Let V be the set of variables.

Input: A finite set of all positive examples in canonically ordering

Output: A finite set of all probabilistic productions

Step 1: Initialize grammar patterns

Two types of productions can be generated by the algorithm are $A \rightarrow \lambda$ and $A \rightarrow aB$

Where A and B are variables and a is a character in the alphabet Σ and λ be an empty string.

Step 2: Initialize the probability of the λ -production

Compute the probability of λ -production (string with length = 0) from the examples.

Step 3: Determine the probability of string with length = $n - 1$ (start from $n = 1$)

Construct production for generating strings with length = n

Consistency checking

Adjust the probability of string with length = $n - 1$.

Step 4: Increase n by 1

Repeat Step 3 again until the end of examples.

Concept of the algorithm: the algorithm starts by computing the probability of all productions corresponding to generate all strings of length n (where n starts from 0 and increasing by 1). For the first step, only one variable S is introduced. For each step ($n > 1$), introduce a new variable if the probability of all productions does not satisfy the consistency of probability constraint, and also adjust the probability values. The PFSM can be constructed using the final productions obtained from the algorithm.

4. Experimental Results

In this section, we will give you an example. Using the set of all positive examples described in Example 8, we will show you how the productions can be obtained from our algorithm.

Step 1: Initialize grammar patterns

Generate two productions: $S \rightarrow \lambda$ and $S \rightarrow aS$. (The variable S states for the starting variable.)

Step 2: Initialize the probability of the λ -production

Compute the probability of λ -production from the example, $\text{prob}(S \rightarrow \lambda) = 0.49$

Step 3: Determine the probability of string with length = 0

Construct the production $S \rightarrow \lambda$ with $\text{prob}(S \rightarrow \lambda) = 0.49 = \text{prob}(\lambda)$.

Step 4: Increase n by 1

Consider strings a and b .

Step 3: Determine the probability of string with length = 1

Construct the production $S \rightarrow aS$ with $\text{prob}(S \rightarrow aS) = 0.2612$

Since $\text{prob}(a) = \text{prob}(S \rightarrow aS) \times \text{prob}(S \rightarrow \lambda)$, then

$$\text{prob}(S \rightarrow aS) = 0.128/0.49 = 0.2612$$

Construct the production $S \rightarrow bS$ with $\text{prob}(S \rightarrow bS) = 0.3469$

Since $\text{prob}(b) = \text{prob}(S \rightarrow bS) \times \text{prob}(S \rightarrow \lambda)$,

$$\text{prob}(S \rightarrow bS) = 0.17/0.49 = 0.3469$$

Using the consistency checking, since the $\text{prob}(S \rightarrow \lambda) + \text{prob}(S \rightarrow aS) + \text{prob}(S \rightarrow bS) \neq 1$,

then the production $S \rightarrow bS$ must be replaced by $S \rightarrow bT$ and $T \rightarrow \lambda$, where $\text{prob}(S \rightarrow bT) = 0.2488$

Since $\text{prob}(b) = \text{prob}(S \rightarrow bT) \times \text{prob}(T \rightarrow \lambda)$, then

$$\text{prob}(T \rightarrow \lambda) = 0.17/0.2488 = 0.6832$$

Step 4: Increase n by 1

Consider strings aa , ab , ba and bb .

Repeat until we obtain the final solutions as follows:

$$\text{prob}(S \rightarrow \lambda) = 0.49 \quad \text{prob}(S \rightarrow aS) = 0.2612 \quad \text{prob}(S \rightarrow bT) = 0.2488$$

$$\text{prob}(T \rightarrow \lambda) = 0.6462 \quad \text{prob}(T \rightarrow aT) = 0.239 \quad \text{prob}(T \rightarrow bS) = 0.1148$$

The PFSM can be constructed as illustrated by Fig. 7.

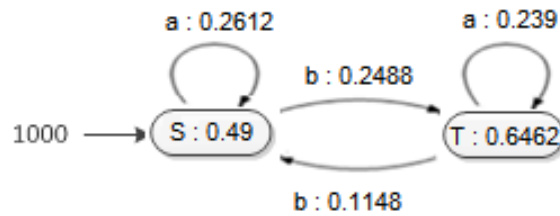


Fig. 7: Probabilistic finite state machine generated by our incremental algorithm.

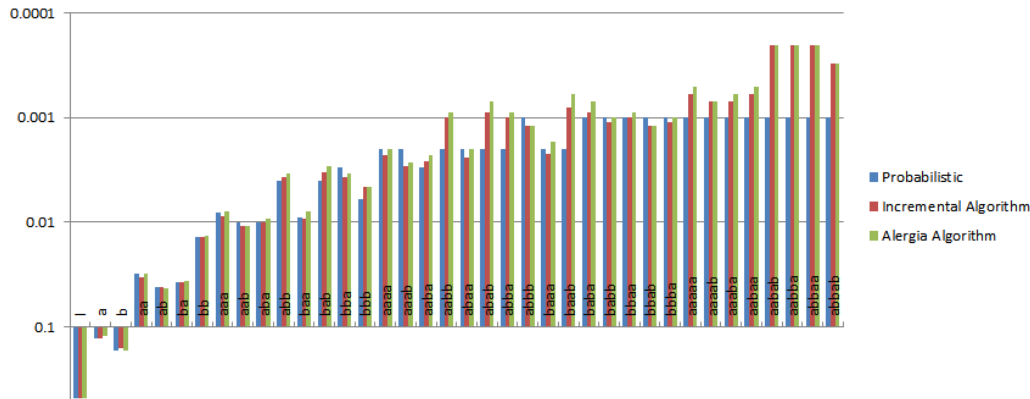


Fig. 8: A comparative bar graph of probabilities between Alergia and our algorithm.

From this comparison graph, we compare probabilities between our incremental algorithm and Alergia algorithm. It is obtained that 20 input examples that their probabilities obtained from our algorithm are closer to the original than those obtained from the Alergia algorithm, 10 input examples that the probabilities obtained from Alergia algorithm are closer than those obtained from our algorithm, and 8 input examples that both algorithms got the same probabilities.

5. Conclusion

We propose an alternative way to construct a probabilistic finite state machine for describing the set of some positive examples of the target language. Our concept is to construct the target probabilistic productions which can fit to the given positive examples. The PFSM can be finally obtained from the obtained productions. We suggest that our algorithm will be an alternative way for constructing a PFSM.

6. References

- [1] M.E. Gold. Language identification in the limit. *Information and Control*. 1967, 10 (5): 447-474.
- [2] N. Chomsky. *Syntactic Structures*. Berlin: Mouton de Gruyter 1957.
- [3] M.V. Zaanen. *Bootstrapping syntax and recursion using alignment-based learning*. Proceedings of the 17th International Conference on Machine Learning, 2000: 1063-1070.
- [4] D. Klein and C.D. Manning. Natural language grammar induction with a constituent-context model. *Pattern Recognition*. 2005, 38 (9): 1407-1419.
- [5] D. Klein. *The unsupervised learning of natural language structure*, Ph.D. thesis, Stanford University, 2005.
- [6] R.C. Carrasco, J. Oncina. Learning stochastic regular grammars by means of a state merging method. In: R.C. Carrasco, J. Oncina (eds.). *Grammatical inference and applications*, Springer, Berlin, Heidelberg 1994: pp. 139-152.
- [7] C. de la Higuera. Learning probabilistic finite automata. In: C. de la Higuera (eds.). *Grammatical Inference: learning automata and grammars*. Cambridge: The United Kingdom at the University Press, 2010: pp. 331-354.
- [8] W. Hoeffding. Probability inequalities for sums of bounded random variables. *American Statistical Association Journal* 1963, 58: 13-30.
- [9] X. Qi. *Analysis on Alergia Algorithm: pattern recognition by automata theory*. Master's Projects, 2016: 491.