

COMPUTATIONAL MODELING OF LANGUAGE ACQUISITION

Rishabh Nigam

Shubhdeep Kochhar

Advisor: Dr. Amitabh Mukherjee

Dept. of Computer Science and Engineering

{rishabh,shubkoch, amit} @ iitk.ac.in

April 14, 2012

ABSTRACT

This project deals with unsupervised learning of Natural Languages. Through a Corpus of sentences which are realistic and natural, we are able to extract pattern out of them and in turn generate new sentences that were not part of original corpus. This Extraction and Generation process is applied on various Corpus of English and Hindi Language and the results are analysed. The Algorithm used for this process is called ADIOS(Automatic Distillation OF Structures). Given a corpus of strings (text or speech, DNA sequencing etc) this algorithm recursively distills a heirarchical structured patterns. For an example if we have sentences like

राम घर जाता है

सीता विद्यालय जाती है

राम विद्यालय जाता है

सीता घर जाती है

from these four sentences the algorithm is able to extract clusters like {राम , सीता} and {घर , विद्यालय} and able to determine pattern such as

{राम , सीता} – {घर , विद्यालय} – {जाता , जाती} – है

INTRODUCTION

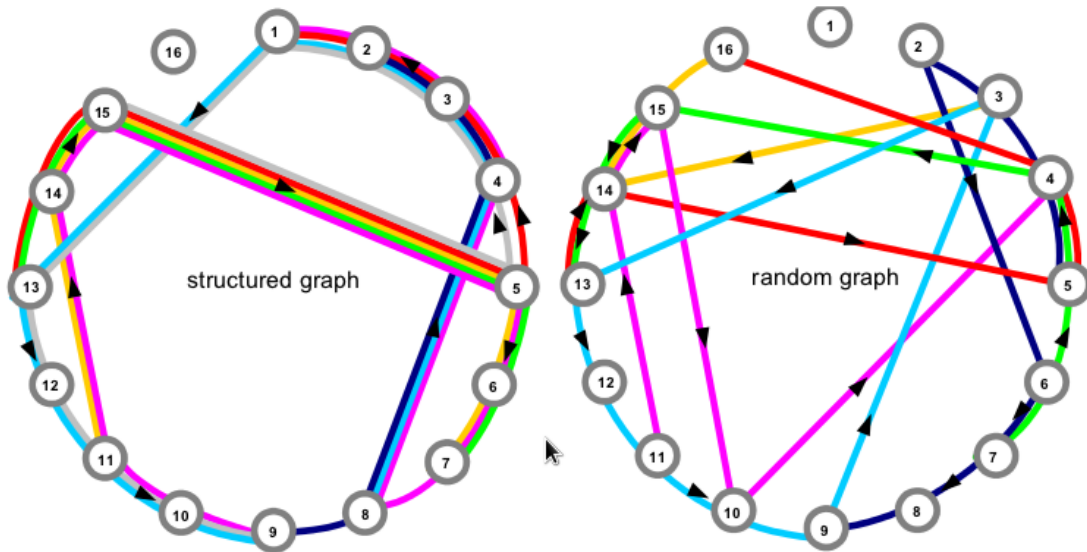
It is well known that the patterns that govern language production are well-formed and rule-like but it is less clear that how are they acquired and what form should such rules take. So in this project, we attempt to learn these rules, and construct patterns in an unsupervised manner. We use probabilistic inference of pattern significance to clusterize similar lexicons and use recursive construction to generate more complex patterns. This way we are able to extract the patterns and rules hidden in the initial corpus of sentences and use these rules and patterns to generate more sentences.

THE ADIOS ALGORITHM

The algorithm which is used in this project is ADIOS(Automatic Distillation Of Structure). The following is a small description of the MEX criterion which is used as a distillation tool for extracting the most significant patterns in the data .

The MEX Criterion

In this criterion we create a pseudograph(a non-simple graph in which both loops and multiple edges are permitted) where the sentences correspond to the different paths and vertices are the unique lexicon entries, augmented by two special symbols 'begin' and 'end'. The following diagram indicates structure that we seek, namely, the bundling of paths, signifying a relatively high probability associated with a sub-structure that can be identified as a pattern.



Reference : Zach Solan Thesis, Tel Aviv University, 2006

Firstly we define a search path $S(e_1 e_2 \dots e_k) = (e_1 ; e_k)$ Next we define two probability functions PR(right) and PL(left) as

$$PR(e_i ; e_j) = p(e_j - e_i e_{i+1} e_{i+2} \dots e_{j-1}) = l(e_i ; e_j) / l(e_i ; e_{j-1})$$

where $l(e_i ; e_j)$ is the number of occurrences of sub-paths $(e_i ; e_j)$ in the graph .

Similarly proceeding from left we have

$$PL(e_j ; e_i) = p(e_i - e_{i+1} e_{i+2} \dots e_{j-1} e_j) = l(e_j ; e_i) / l(e_j ; e_{i+1})$$

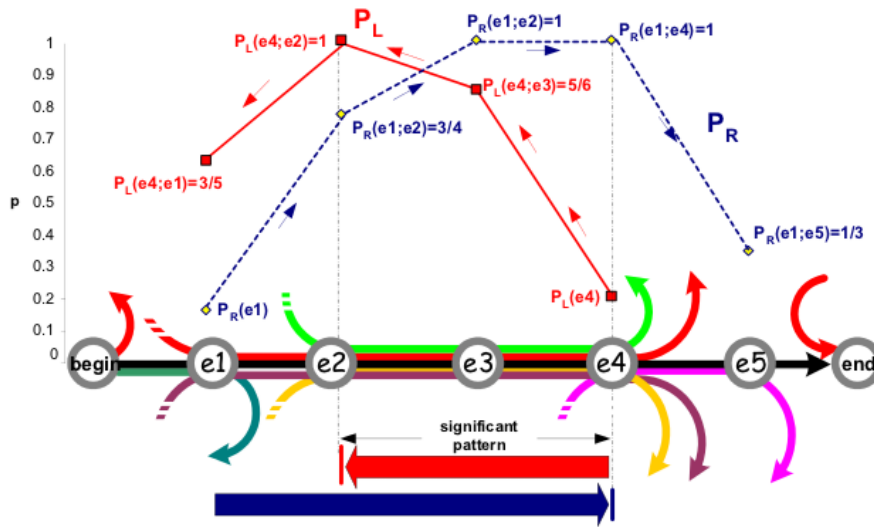
We define decrease ratio, DR $(e_i ; e_j)$, whose value at e_j is $DR(e_i ; e_j) = PR(e_i ; e_j) / PR(e_i ; e_{j-1})$ and another decrease ratio, DL $(e_j ; e_i) = PL(e_j ; e_i) / PL(e_{j+1} ; e_i)$

The algorithm calculates PL and PR from all possible starting points and this defines a matrix of the form

$$M = PR(e_i ; e_j) \quad i > j$$

$$PL(e_j ; e_i) \quad i < j$$

$$P(e_i) \quad i = j$$



Reference : Zach Solan Thesis, Tel Aviv University, 2006

	<i>begin</i>	e_1	e_2	e_3	e_4	e_5	<i>end</i>
<i>begin</i>	8/41	2/4	1/3	1/3	1/3	1	1
e_1	2/8	4/41	3/6	3/5	3/5	1	1
e_2	1/2	3/4	6/41	1	1	1	1
e_3	1	1	5/6	5/41	5/6	1/2	1
e_4	1	1	1	1	6/41	1	1/2
e_5	1	1/3	1/5	1/5	2/6	2/41	1
<i>end</i>	1	1	1	1	1/2	1	8/41

Reference : Zach Solan Thesis, Tel Aviv University, 2006

The algorithm then identifies the leading pattern from the matrix by examining the DR and DL and is returned as the outcome of search in question.

The ADIOS Algorithm

The algorithm works in three steps:

1. **INITIALIZATION**: sentence loading
2. **PATTERN DISTILLATION**: iterative search for significant patterns which are added to lexicons as new units
3. **GENERALIZATION**: it generates more and more candidate pattern to be considered by pattern distillation

So an overview of ADIOS Algorithm is as follows:

1. **Initialization** (load all sentences)
2. repeat
3. for all $m = 1 : N$ do N is the number of paths in the graph
 - Pattern Distillation(m)** (Identifies new significant patterns in search-path m using the MEX criterion)
 - Generalization(m)** (Generate new pattern candidates for search-path (m))end for
4. until no further significant patterns are found

Reference : Zach Solan Thesis, Tel Aviv University, 2006

DATABASE USED

We tried the algorithm on various corpus :

- First one was a corpus of 500 sample sentences that was provided along with the codes. This showed great results as it had sentences with similar structures. It also contained a context-free-grammar file (.cfg) which is useful in calculating the precision and recall , the factors that are used to assess the performance of the algorithm.
- Secondly we tried the code on the CHILDES corpus. The corpus contained the conversation of a mother and a child who is in the early stages of language aquisition. We chose this database to get a better understanding of natural language aquisition in a child. From the database we extracted about 24,000 sentences and generated grammar from it.
- Then, we moved on to HINDI database which we got from CFILT (Center For Indian Language Technology,IIT Bombay). We used two different sets of sentences, from one we used 2,500 sentences and from other we used 5,000 sentences to generate the grammar .
- Lastly we used another small database of a commentary which we recieved from Mr. Shushobhan which consisted of about 800 sentences. The database consisted of commentary describing a certain set of processes and the sentences were quite similar in structure and grammar.

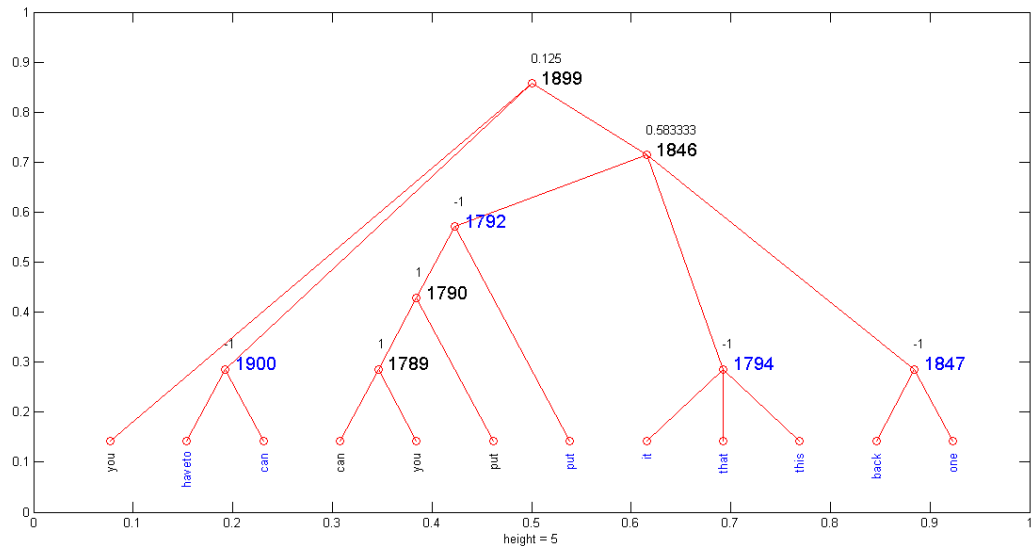
IMPLEMENTATION DETAILS

The first part consisted of preprocessing of the data used, for this we had to write a script to extract the Mother Sentences from the CHILDES database. The CHILDES database consisted of sentences of Mother as well as some information related to the sentences, so we had to extract only the sentences used by the mother. Then as we had most of the code for the ADIOS algorithm, we basically had to understand and use it. For this we really had to have an understanding of the algorithm as the code had certain parameters which were initially not very clear. Then as the code had lots of commands, we needed to decide on what to do first (in this the documentation http://adios.tau.ac.il/algorithm.html#Step_By_Step was of great help). Then, we needed to print the patterns obtained after the completion of the algorithm for which we used some matlab codes (these codes were made available from the site <http://adios.tau.ac.il/algorithm.html>). Then we decided to run this on a Hindi database. We got the database from CFILT IIT Mumbai. Now this database again needed to be converted into suitable format, for this again we needed to do some manipulations with the file that we got. After the preprocessing stage we needed to train the grammar on the database and generate new sentences. For this we again had to repeat the steps we used for the CHILDES database. We decided to use another set of sentences, this time doubled the number of sentences from 2500 to 5000. We also used the algorithm on a Commentary of around 780 sentences.

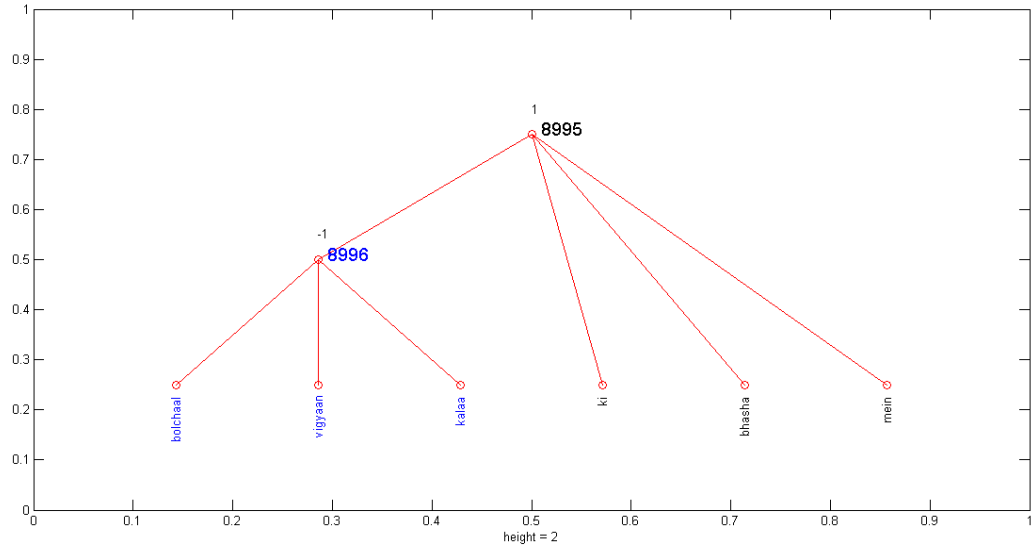
RESULTS

For the CHILDES database, we were able to get 101 equivalence classes. For the Hindi database we were able to generate 10 equivalence classes, for the commentary we got 14 equivalence classes. Some of the equivalence classes are shown in the pictures below, the complete result can be found at home.iitk.ac.in/~rishabh/cs365/projects/{hindi,childes,commentary}.{label,generate}.txt

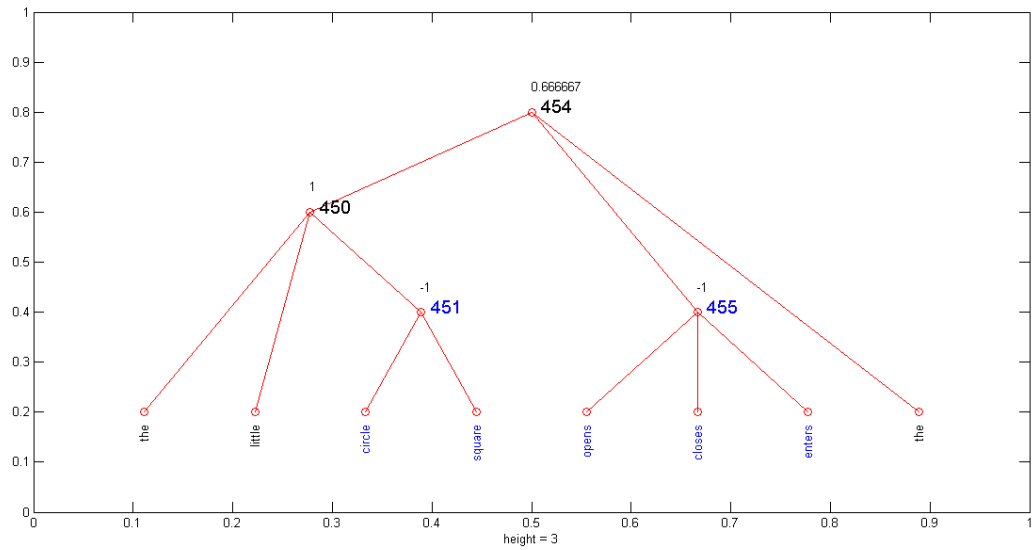
CHILDES



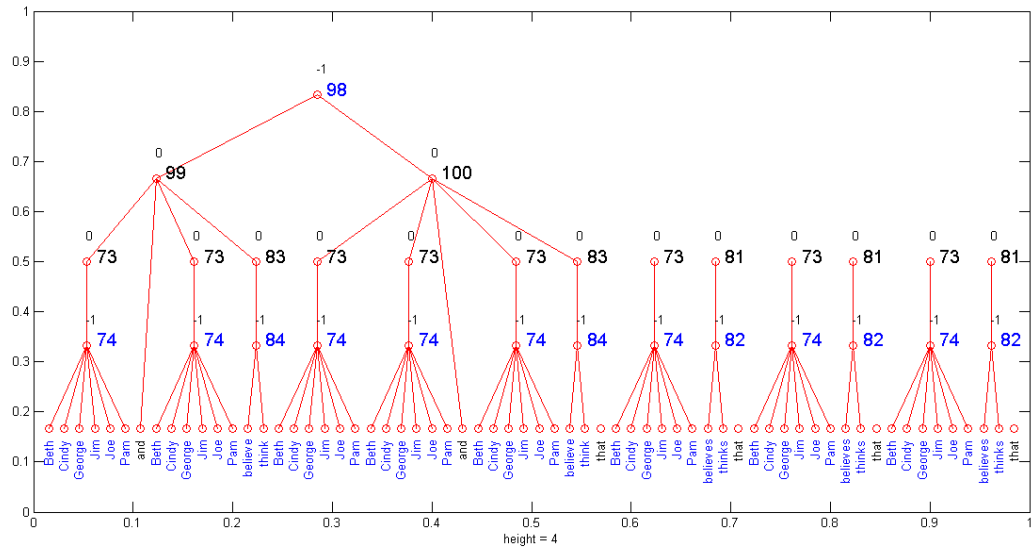
HINDI



COMMENTARY



SAMPLE



BIBLIOGRAPHY

- [1] Heider. Waterfall ,Ben Sandbank, Luca Onnis and Shimon Edelman , An empirical generative framework for computational modeling of language acquisition : Cambridge University Press 2010 <http://kybele.psych.cornell.edu/~edelman/Waterfall-Sandbank-Onnis-Edelman-JCL10.pdf>

- [2] Zach Solan PHD thesis, Unsupervised Learning of Natural Languages, under Professor David Horn, Professor Shimon Edelman, Professor Eytan Ruppim : Senate of Tel Aviv University 2006 [1-38] <http://horn.tau.ac.il/publications/ZachSolanThesis.pdf>

- [3] CHILDES (Child Language Data Exchange System) <http://childes.psy.cmu.edu/>

- [4] Dr Pushpak Bhattacharya, CFILT(Centre For Indian Language technology) IIT Bombay <http://www.cfilt.iitb.ac.in/Downloads.html>

- [5] ADIOS website by Zach Solan, Tel Aviv University <http://adios.tau.ac.il/algorithm.html>