

# Incremental Discovery of Sequential Pattern from Semi-structured Document Using Grammatical Inference

Ramesh Thakur<sup>1</sup>, Suresh Jain<sup>2</sup>, and Narendra S. Chaudhari<sup>3</sup>

<sup>1</sup> IIPS, DAVV, Indore, India

<sup>2</sup> KCB Technical Academy, Indore, India

<sup>3</sup> Indian Institute Of Technology, Indore

**Abstract.** On the World Wide Web a large numbers of information is available in the form of semi-structured format. Knowledge discovery in semi-structured document has been recognized as promising task. Since semi structured document is typically hidden within HTML formatting intended for human viewing the details of which vary widely from site to site and frequent changes made to their formatting so we can't construct a global schema, discovery of interesting rules from it is complex and tedious process. Most of the existing system uses hand-coded wrappers to extract information, which is monotonous and time consuming. An intelligent and automated method is needed for their processing. Learning grammatical information from given sample of semi-structured documents has attracted lots of attention in the past decades. To understand "what say the data" is necessary to know the structure of data to discover the syntactic-semantic knowledge of its language.

The problem of learning the correct grammar for the unknown language from finite example of the language is known as grammatical inference problem. In automated grammar learning, the task is to infer grammar rules from given information about the target language. If example belongs to the target language it is called positive example otherwise it is called negative example. In this paper we propose a grammar inference methodology to automate the construction of grammar rules and facilitate the process of information extraction. We are using hybrid technique of association analysis and sequential algorithm to generate context free grammar rules from semi-structured document (HTML document).

Our algorithm that infers a sequential pattern from a sequence<sup>1</sup> of discrete HTML tags. The basic insight is that sub-string is selected on the basis of high support factor<sup>2</sup> by taking entire sentences into account. Which appears more frequently in string can be replaced by a grammatical rule that generate the sub-string, and this process is repeated many times, producing a single length rules of the sequence. The result is strictly a context-free grammar rules, which provide a compact summary of corpora that aids understanding of its properties.

**Keywords:** Knowledge discovery, sequence mining and grammar inference.

---

<sup>1</sup> Sequence: A sequence is an ordered list of alphabet symbols. We denote a sequence by  $\langle s_1s_2\dots s_n \rangle$  where  $s_i$  is a symbol. A sequence  $\langle a_1, a_2, \dots, a_n \rangle$  is a subsequence of another sequence  $\langle b_1, b_2, \dots, b_m \rangle$  if there exist integers  $i_1 < i_2 < \dots < i_n$  such that  $a_1 = b_{i_1}, a_2 = b_{i_2}, \dots, a_n = b_{i_n}$ . For example, the sequence bob is subsequence of bbobbb.

<sup>2</sup> Support factor: The support Factor (SF $\beta$ ) for sub-sequences in corpora C.

$$SF\beta = \frac{\sum_{i=1}^N \text{count of } \beta \text{ in sentence} \times \text{length of } \beta}{\text{length of sentence}}$$

Where N= number of sentences in Corpora C and  $\beta$  is a candidate sub-sequence for replacement.