*Semantic Web Project Report on*

# Autosummarization of Technical Paper using Grammatical Rule Inference

**S.Gopinath (12IT64)**

**S.Vinay Kumar(12IT66)**

Under the Guidance of,

**Mr. Dinesh Naik**

Department of Information Technology, NITK Surathkal

*Date of Submission: 8th May, 2015*

in partial fulfillment for the award of the degree of

**Bachelor of Technology**

In

**Information Technology** At



**Department of Information Technology**

**National Institute of Technology Karnataka, Surathkal**

**April 2015**

# Department of Information Technology, NITK Surathkal
## Semantic Web Technology Project
## End Semester Evaluation Report (April 2015)

---

**Course Code :** IT 412

**Course Title:** Semantic Web Technology

**Project Title:** *Autosummarization of Technical Paper using Grammatical Rule Inference*

**Project Group:**

| Name of the Student | Register No. | Signature with Date |
|---|---|---|
| S.Gopinath | 12IT64 | |
| S.Vinay Kumar | 12IT66 | |

Place:

Date:                                                          *(Name and Signature of Guide)*

# Abstract

Over the years, the amount of information available electronically has grown manifold. There is an increasing demand for automatic methods for text summarization. Domain independent techniques for automatic summarization by paragraph extraction have been proposed in many papers [3, 4]. In this paper, we attempt to summarize a technical paper written in html using grammatical rule inference technique so that it becomes easier for us humans to compare between many technical papers just by reading their summary.

**Keywords:** *Auto summarization, Grammar Rule Inference Technique, Semantic Web*

# Contents

# List of Figures

# 1    Introduction

Nowadays the Web has proved to be as a rich and extraordinary data source of information, here multiple domains can be accessed and mined. Mining Web data is referred as Web Mining. Some of the objectives of mining web data include finding relevant information discovering new knowledge from web personalized, web synthesis and learning about individual users. Amongst these the most common use is finding relevant information. We simply specify a set of keywords or query as a request or a reference and we get a list of pages, ranked as per similarity of query.

Currently searching web face with one problem that many times outcome is not satisfactory because of irrelevance of the information. Searching the exact information from such a huge repository [2] of unstructured web data is still main area of research interest. One solution to this problem is Semantic Web. The Semantic Web is an extension of current Web in which information is given as well defined meaning, hence enabling computers and people to work with better coordination. By using the existing web semantically new semantically structure can be exploited, and then the results of web mining can be improved, thereby building semantic Web [1].

Automatic summarization is the process of reducing a text document with a computer program in order to create a summary that retains the most important points of the original document. As the problem of information overload has grown, and as the quantity of data has increased, so has interest in automatic summarization. Technologies that can make a coherent summary take into account variables such as length, writing style and syntax

Generally, there are two approaches to automatic summarization: extraction and abstraction. Extractive methods work by selecting a subset of existing words, phrases, or sentences in the original text to form the summary. In contrast, abstractive methods build an internal semantic representation and then use natural language generation techniques to create a summary that is closer to what a human might generate. Such a summary might contain words not explicitly present in the original. Research into abstractive methods is an increasingly important and active research area, however due to complexity constraints, research to date has focused primarily on extractive methods.

## 1.1  *Motivation*

Current automatic summarizers usually rely on sentence extraction to produce summaries. Human professionals also often reuse the input documents to generate summaries; however, rather than simply extracting sentences and stringing them together, as most current summarizers do, humans often "edit" the extracted sentences in some way so that the resulting summary is concise and coherent.

# 2     Literature Survey

Currently works have been done in two common approaches to summarisation. The first approach comprises of analysis of text, and rewriting or rephrasing it in a short way. Until today this approach hasn't achieved any substantial results. The second approach, which is similar to algorithm used in this project, tries to extract the key sentences from the text, and then tries to put them together properly. One famous algorithm that implements this approach is TextRank [5].

TextRank is a general purpose graph-based ranking algorithm for NLP. Essentially, it runs PageRank on a graph specially designed for a particular NLP task. For keyphrase extraction, it builds a graph using some set of text units as vertices. Edges are based on some measure of semantic or lexical similarity between the text unit vertices. Unlike PageRank, the edges are typically undirected and can be weighted to reflect a degree of similarity. Once the graph is constructed, it is used to form a stochastic matrix, combined with a damping factor (as in the "random surfer model"), and the ranking over vertices is obtained by finding the eigenvector corresponding to eigenvalue 1 (i.e., the stationary distribution of the random walk on the graph).

There also have been works where one can perform a robust decomposition of the text of an article into individual segments that correspond to basic natural parts, such as title, section title, text, etc. This was aimed at automation of the task of building an Ontology of documents. We combine these two approaches for summarisation of a technical papers.

## 2.1   *Problem Statement*

"To develop a Travel Route Recommendation System for Indian cities, incorporating user friendly interaction through natural speech dialogues, and natural language processing techniques"

# 3    Methodology

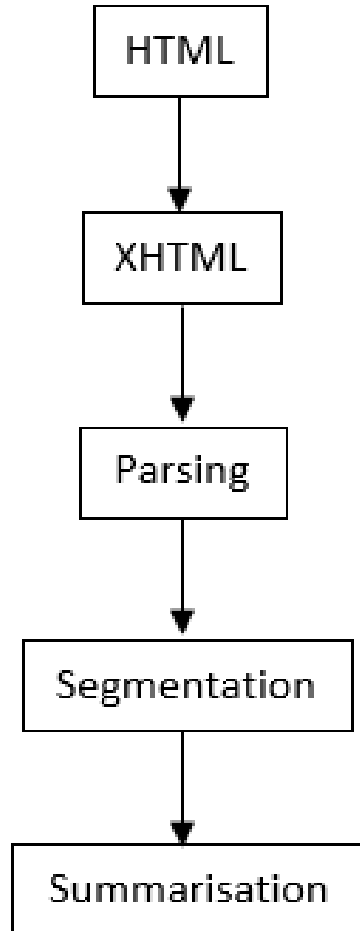## 3.1   *System Architecture*



Figure 3.1: System Architechture

## 3.2   *Different Components of the Project*

1. Preprocessing:

   The extraction of information from web pages is the subject matter of the Wrapper Induction work, although their purpose is not the inference of an ontology, but to enable a database like use of that information. In particular, HTML aware wrapper induction approaches bear similarities to the first step of the method developed In particular, they share the preprocessing of HTML pages with the objective of

4

cleaning frequent syntax errors and transformation into valid XML syntax (using the XHTML format) to ease further processing.

2. Summarization:

**Intersection function:** This function receives two sentences, and returns a score for the intersection between them. We just split each sentence into words/tokens, count how many common tokens we have, and then we normalize the result with the average length of the two sentences.

**The sentences dictionary:** This part is actually the Heart of the algorithm. It receives our text as input, and calculates a score for each sentence. The calculations is composed of two steps: In the first step we split the text into sentences, and store the intersection value between each two sentences in a matrix (two-dimensional array). So values[0][2] will hold the intersection score between sentence 1 and sentence 3.

In fact, we just converted our text into a fully-connected weighted graph! Each sentence is a node in the graph and the two-dimensional array holds the weight of each edge!

In the second step we calculate an individual score for each sentence and store it in a key-value dictionary, where the sentence itself is the key and the value is the total score. We do that just by summing up all its intersections with the other sentences in the text (not including itself).

We calculates the score for each node in our graph. We simply do that by summing all the edges that are connected to the node. Building the summary: Obviously, the final step of our algorithm is generating the final summary. We do that by splitting our text intoparagraphs, and then we choose the best sentence from each paragraph according to our sentences dictionary. The Idea here is that every paragraph in the text represents some logical subset of our graph, so we just pick the most valuable node from each subset!

Why this works: There are two main reasons why this algorithm works: The first reason is that a paragraph is a logical atomic unit of the text. In simple words there is probably a very good reason why the author decided to split his text that way. The second reason is that if two sentences have a good intersection, they probably holds the same information.

So if one sentence has a good intersection with many other sentences, it probably holds some information from each one of them- or in other words, this is probably a key sentence in our text.

# 4    Work Done

A technical paper in HTML is taken as the input. The HTML code may be having inconsistencies which may lead to problems. Hence the HTML document is standardized by converting it into a XHTML document. This makes the document structured and it get easier for extraction of segments. After converting the HTML document into XHTML we segment the document and extract paragraphs with the help of ¡p¿¡/p¿ tags that define a paragraph. Once we get the segments (paragraphs), we use our summarization algorithm to summarize each of these paragraphs. We combine the summaries of each paragraphs and create a new document containing the summaries.

# 5    Results and Analysis

A sample input HTML paper that has been taken as input is shown in the Figure 5.1, Figure 5.2 and Figure 5.3.



Figure 5.1: Input Data-HTML Page-Part 1



Figure 5.2: Input Data-HTML Page-Part 2

The output, i.e. the summary of the given HTML is shown Figure 5.4 and Figure 5.5.

Figure 5.3: Input Data-HTML Page-Part 3



Figure 5.4: Output - Summary of the Technical Paper - Part 1

# 6 Conclusion & Future Work

An application has been created using which a user can get the summary of any technical paper written in HTML so that it becomes easy for the user to compare between many such papers for similarities.

## 6.1 *Future Work*

Future work includes functionalities that can compare the summaries of any two technical papers automatically and tell the user whether those two papers are similar or not.

Figure 5.5: Output - Summary of the Technical Paper - Part 2

# References

[1] Dana Angluin. Inference of reversible languages. *Journal of the ACM (JACM)*, 29(3):741–765, 1982.

[2] Valter Crescenzi, Giansalvatore Mecca, Paolo Merialdo, et al. Roadrunner: Towards automatic data extraction from large web sites. In *VLDB*, volume 1, pages 109–118, 2001.

[3] Jade Goldstein, Vibhu Mittal, Jaime Carbonell, and Mark Kantrowitz. Multi-document summarization by sentence extraction. In *Proceedings of the 2000 NAACL-ANLPWorkshop on Automatic summarization-Volume 4*, pages 40–48. Association for Computational Linguistics, 2000.

[4] Mandar Mitray, Amit Singhalz, and Chris Buckleyyy. Automatic text summarization by paragraph extraction. *Compare*, 22215(22215):26, 1997.

[5] Wikipedia. Automatic summarization — wikipedia, the free encyclopedia, 2015. [Online; accessed 8-May-2015].