



A Review Paper: Categorization of Web Pages

Ritu¹, Sapna Aggarwal²

Student, M.TECH(CSE) Jind Institute of Engineering & Technology, Jind, India¹

Assistant Professor, (CSE) Jind Institute of Engineering and Technology, Jind, India²

mail_id@gmail.com¹

er.sapna.aggarwal@gmail.com²

Abstract

Contemporary web is comprised of trillions of pages and everyday tremendous amount of requests are made to put more web pages on the WWW. It has been difficult to manage information present on web than to create it. Web page categorization can be defined as an approach to categorize the web pages based on a set of predefined categories to manage large web content. Yahoo! and ODP are the examples of web directories in which pages are categorized manually or semiautomatically, but it is a very time consuming task. There are many ways of categorizing web pages using different techniques. An approach to categorize web pages automatically on the basis of characteristics of web pages using neural network based single discrete perceptron training algorithm which is extended by selecting webpage specific features to categorize web pages of predefined categories with high accuracy. The idea is presented with the help of two specific and major categories of web pages chosen for categorization that are newspaper and education.

Keywords: *Categorization, Web*

1. Introduction

Internet is the source of enormous amount of information accessed by large number of people every day. Contemporary web is comprised of trillions of pages and everyday tremendous amount of requests are made to put more web pages on the WWW. It has been difficult to manage information present on web than to create it. Web page categorization can be defined as an approach to categorize the web pages based on a set of predefined categories to manage large web content. Yahoo! and ODP are the examples of web directories in which pages are categorized manually or semiautomatically, but it is a very time consuming task. There are many ways of categorizing web pages using different techniques. An approach to categorize web pages automatically on the basis of characteristics of web pages using neural network based single discrete perceptron training algorithm which is extended by selecting webpage specific features to categorize web pages of predefined categories with high accuracy.

The growing number of applications on the web leads to rapid increase in number of web pages. The data available on the web can be in the form of

text, images, audio, video, graphics and many other forms. Web pages present on the web can be static or dynamic. The content of dynamic web pages keeps on changing time to time. Web is considered as a large repository of information which is accessed by millions of users' everyday through internet. The dynamic nature of web and large scale explosion of web pages may put a threat to efficient information retrieval tasks. Web can be considered as an information resource, therefore it is important to describe and organize the huge content present on the web in order to realize web's full potential. Thus web page categorization is an intellectual task, important and indeed essential for organizing and understanding web content for different applications, efficient information retrieval and other tasks related to web mining. Here we will discuss some facts about web page categorization including the types of web page categorization, need of web page categorization and various characteristics of web pages.

2. Literature Review

From the very beginning categorization was done manually by domain experts. Yahoo! [3] and ODP [4] are the examples of web directories which are developed manually. But with the rapid increase of web pages it became extremely difficult to categorize web pages manually. Therefore categorization began to be done semi automatically or automatically. There are a number of approaches which have been applied in the field of web page categorization including K-Nearest Neighbor approach [11], Bayesian probabilistic models [12], inductive rule learning, decision trees, neural networks and support vector machine. All the above mentioned approaches

ISSN : 2348-5612 © URR



9 770234 856124