

## Performance Evaluation of RM-MapReducer in Web Page Categorization

<sup>1</sup>P. Malarvizhi and <sup>2</sup>N. Radhika

<sup>1</sup>Department of CSE, Karpagam University, Coimbatore, Tamil Nadu, India

<sup>2</sup>Department of Computer Science and Engineering, Amrita School of Engineering,  
Amrita University, Coimbatore, Amrita Vishwa Vidyapeetham, India

---

**Abstract:** This study presents the performance evaluation of Relevancy Measure based MapReducer (RM-MapReducer) in supporting web page categorization with different datasets between 200 and 400 web pages as test data. Experiments were performed using datasets of WebKB and ODP with the real time crawler dataset. It was observed that the real time crawler dataset provides better results compared to the other dataset of WebKB and ODP using MapReduce programming model.

**Key words:** Web page categorization, RM-MapReducer, dataset, web pages, programming model

---

### INTRODUCTION

The dimension and the dynamic nature of the web creates a need for classification of web pages (Malarvizhi and Radhika, 2016) to manage the information retrieval process. The rapid growth of web increases the need to have automated assistance to organize the huge amount of information provided by keyword-based search engines. Categorization of web pages is a process that assigns the labels of the category of predefined to a web page (Malarvizhi and Radhika, 2016). The explosive growth of web introduces a need for web page classification (Uysal, 2016) to manage the information retrieval tasks and to improve the performance of the search (Qi and Davison, 2009). In this research, MapReduce parallel programming model with functions map and reduce is used for categorizing the web pages based on the relevancy measure. MapReduce, the Google's popular framework is an attractive parallel model with functions map and reduce suitable for parallel processing of arbitrary data (Malarvizhi and Pujeri, 2012).

### MATERIALS AND METHODS

**Relevancy Measure based MapReducer (RM-MapReducer):** MapReduce, the parallel programming model was used to categorize the web pages based on the relevancy measure. The Relevancy Measure RM is computed by comparing the similarities between the two vectors  $v_1$  of frequent keyword and  $v_2$  of keyword with the predefined category vector  $C$  is given below. For all  $w \in W$  do; for all  $c \in C$  do;  $c \leftarrow (RM)_{\max}$ .

RM is the sum of the two similarities of  $v_1$  and  $v_2$  (Malarvizhi and Pujeri, 2012). MapReduce is a popular programming model with functions map and reduce and meets a number of varieties of applications (Chen and Schlosser, 2008). The MapReduce programming model assigns a predefined category label to a web page based on the (RM) max value. The dataset of web pages are given as input to the map function of the MapReducer and the reduce function of the MapReducer computes the relevancy measure for each category and assigns the (RM)max value category label to the web pages.

### RESULTS AND DISCUSSION

This comparative study was performed on dataset of 400 web pages and was implemented using java. The performance of the study was evaluated on test data with the evaluation metrics of precision and recall.

**Dataset:** Two datasets, dataset 1 and 2 are created for test data. Dataset 1 contains the web pages related to the category of C1-4 of course, student, department and conference. Web pages of category course, student and department were collected from WebKB dataset and category conference was collected from computer science conferences of the ODP website. Dataset 2 contains the web pages of category C1-4 of course, student, department and conference collected from web by using the web crawler websphinx. Websphinx is a customized web crawler used to collect the web pages from web. The 200 web pages were collected for each dataset of C1 of 62 web pages, C2 of 62 web pages, C3 of 35 web pages, C4 of 41 web pages and totally 400 web pages were collected for dataset.