Chapter XII Feature Selection for Web Page Classification

K. Selvakuberan *Tata Consultancy Services, India*

M. Indra Devi *Thiagarajar College of Engineering, India*

R. Rajaram *Thiagarajar College of Engineering, India*

ABSTRACT

The World Wide Web serves as a huge, widely distributed, global information service center for news, advertisements, customer information, financial management, education, government, e-commerce and many others. The Web contains a rich and dynamic collection of hyperlink information. The Web page access and usage information provide rich sources for data mining. Web pages are classified based on the content and/or contextual information embedded in them. As the Web pages contain many irrelevant, infrequent, and stop words that reduce the performance of the classifier, selecting relevant representative features from the Web page is the essential preprocessing step. This provides secured accessing of the required information. The Web page. This information may be used to personalize the information needed for the users and to preserve the privacy of the users by hiding the personal details. The issue lies in selecting the features which represent the Web pages and processing the details of the user needed the details. In this chapter we focus on the feature selection, issues in feature selection, and the most important feature selection techniques described and used by researchers.

INTRODUCTION

There are an estimated 15 to 30 billion pages available in the World Wide Web with millions of pages being added daily. Describing and organizing this vast amount of content is essential for realizing the Web's full potential as an information resource. Automatic classification of Web pages is needed for the following reasons. (a) Large amount of information available in the internet makes it difficult for the human experts to classify them manually (b) The amount of Expertise needed is high (c) Web pages are dynamic and volatile in nature (e) More time and effort are required for classification. (f) Same type of classification scheme may not be applied to all pages (g) More experts needed for classification. Web page classification techniques use concepts from many fields like Information filtering and retrieval, Artificial Intelligence, Text mining, Machine learning techniques and so on. Information filtering and retrieval techniques usually build either a thesauri or indices by analyzing a corpus of already classified texts with specific algorithms. When new text is to be classified, thesaurus and index are used to find the similarity with already existing classification scheme to be associated with this new text.

Until the late 1980s, the most effective approach to Web page classification seemed to be that of manually by building classification systems by means of knowledge-engineering techniques, i.e. manually defining a set of logical rules that encode expert knowledge on how to classify Web page documents under the given set of categories. In the 1990s this perspective has been overturn, and the machine learning paradigm to automated Web page classification has emerged and definitely superseded the knowledge-engineering approach. Within the machine learning paradigm, a general inductive process automatically builds an automatic text classifier by "learning", from a set of previously classified Web documents, the characteristics of the categories of interest. The advantages of this approach are accuracy comparable to human performance and a considerable savings in terms of expert manpower, since no intervention from either knowledge engineers or domain experts is needed. Currently Web page categorization may be seen as the meeting point of machine learning and information retrieval. As Machine Learning aims to address larger, more complex tasks, the problem of focusing on the most relevant information in a potentially overwhelming quantity of data has become increasingly important. For instance, data mining of corporate or scientific records often involves dealing with both many features and many examples, and the internet and World Wide Web have put a huge volume of low-quality information at the easy access of a learning system. Similar issues arise in the personalization of filtering systems for information retrieval, electronic mail, net news, etc.

The main objective of this chapter is to focus on the feature selection techniques, need for feature selection, their issues in Web page classification, feature selection for privacy preserving data mining and the future trends in feature selection.

LITERATURE SURVEY

Rudy Setiono and Huan Liu (1997) proposed that Discretization can turn numeric attributes into discrete ones. χ^2 is a simple algorithm. Principal Component Analysis-compose a small number of new features. It is improved from simple methods such as equi-width and equal frequency intervals. For each and every attributes calculate the χ^2 value for each and every interval. Combine the lowest interval values while approximation.

Shounak Roychowdhury (2001) proposed a technique called granular computing for processing and expressing chunks of information called granules. It reduces hypothesis search space, to reduce storage. Fuzzy set based feature elimination techniques in which subset generation and subset 14 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/feature-selection-web-page-classification/29152

Related Content

Data Mining and the Project Management Environment

Emanuel Camilleri (2010). *Data Mining in Public and Private Sectors: Organizational and Government Applications (pp. 337-357).* www.irma-international.org/chapter/data-mining-project-management-environment/44296

Analysis and Integration of Biological Data: A Data Mining Approach using Neural Networks

Diego Milone, Georgina Stegmayer, Matías Gerard, Laura Kamenetzky, Mariana Lópezand Fernando Carrari (2013). *Data Mining: Concepts, Methodologies, Tools, and Applications (pp. 203-230).* www.irma-international.org/chapter/analysis-integration-biological-data/73441

Conceptual Clustering of Textual Documents and Some Insights for Knowledge Discovery

Leandro Krug Wives, José Palazzo Moreira de Oliveiraand Stanley Loh (2008). *Emerging Technologies of Text Mining: Techniques and Applications (pp. 223-243).* www.irma-international.org/chapter/conceptual-clustering-textual-documents-some/10184

An Envisioned Approach for Modeling and Supporting User-Centric Query Activities on Data Warehouses

Marie-Aude Aufaure, Alfredo Cuzzocrea, Cécile Favre, Patrick Marceland Rokia Missaoui (2013). International Journal of Data Warehousing and Mining (pp. 89-109). www.irma-international.org/article/envisioned-approach-modeling-supporting-user/78288

Analysis on the Allocation of Control Right of Intergenerational Inheritance of Family Enterprises in the New Era

Zhanzhong Wangand Jiajun Li (2023). *International Journal of Data Warehousing and Mining (pp. 1-20).* www.irma-international.org/article/analysis-on-the-allocation-of-control-right-of-intergenerational-inheritance-of-familyenterprises-in-the-new-era/319966