

Chapter 7

Document Classification

ABSTRACT

Keywords can be used as attributes for mining rules or as a basis for measuring the similarity of new (unclassified) documents with existing (classified) ones. The focus is on the problem of extracting keywords from document collection in order to use them as attributes for document classification. Document classification is a hot topic in machine learning. Typical approaches extract “features,” generally words, from document, and use the feature vectors as input to a machine learning scheme that learns how to classify documents. This “bag of keywords” model neglects keyword order and contextual effects.

DOCUMENT CLASSIFICATION: ROLE OF KEYWORDS

Survey shows that text classification is a typical scholarly activity in literary study, and automatic text classification methods can be used in three scenarios.

1. The first is **Information Organization** - A classifier can learn the target category concepts (e.g. news article about trade, acquisition, etc.) from the training documents, and then assign new documents into these predefined categories.
2. The second purpose is **Knowledge Discovery** - A successful classifier can provide insights to understand a target concept by revealing the correlations between the features and the concept.

DOI: 10.4018/978-1-7998-3772-5.ch007

Document Classification

3. The third purpose is **Example-based Retrieval** - A classifier might be able to learn a concept from a small number of training documents with the help of semi-supervised learning or active learning methods, and then retrieve more documents similar to the training examples from a large collection.

Automatic document classification methods are tremendously used in research domain of textual documents due to,

- The numerous and important domains of application
- The indispensability in many applications
- The implausibility of manual alternative
- The productivity of machine learning
- Large amount of digital form of documents

NEED FOR DOCUMENT CLASSIFICATION

The goal of text classification is to (semi) automate the categorization process. It is also useful for reducing cost and improving performance (including accuracy and consistency) of text processing.

Scope of Text Classification

Text Classification could be considered as the application of (semi) automatic methods in order to choose, from a set of predefined classification codes, the appropriate one (category / class) for a given new document. Various studies have focused on the construction of a model (Rule or Tree) related to the existence of keywords in order to assign the class of the unclassified document. The text collection is divided into documents and each document is characterized by a number of keywords (Yang, Y., and Pedersen, J.O, 1997). In automatic document classification, for example, for classifying newspaper articles into predefined categories such as politics and sports, the crucial step is how to select appropriate keywords. With traditional classification methods based on the vector space model, frequent words are emphasized and therefore low- frequency words tend to be disregarded. However, there often exist low-frequency words that are effective for classification. For instance,

3 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/document-classification/268466

Related Content

The Power of Sampling and Stacking for the PaKDD-2007 Cross-Selling Problem

Paulo J.L. Adeodato, Germano C. Vasconcelos, Adrian L. Arnaud, Rodrigo C.L.V. Cunha, Domingos S.M.P. Monteiro and Rosalvo F.O. Neto (2008). *International Journal of Data Warehousing and Mining* (pp. 22-31).

www.irma-international.org/article/power-sampling-stacking-pakdd-2007/1804

The Use of Smart Tokens in Cleaning Integrated Warehouse Data

Christie I. Ezeife and Timothy E. Ohanekwu (2005). *International Journal of Data Warehousing and Mining* (pp. 1-22).

www.irma-international.org/article/use-smart-tokens-cleaning-integrated/1749

Retailer Case Study

Johanna Wenny Rahayu, David Tanier and Eric Pardede (2006). *Object-Oriented Oracle* (pp. 276-323).

www.irma-international.org/chapter/retailer-case-study/27343

Mining Organizations' Networks: Multi-Level Approach1

James A. Danowski (2012). *Social Network Mining, Analysis, and Research Trends: Techniques and Applications* (pp. 205-230).

www.irma-international.org/chapter/mining-organizations-networks/61520

Application of Artificial Neural Network and Genetic Programming in Civil Engineering

Pijush Samui, Dhruvan Choubisa and Akash Sharda (2014). *Biologically-Inspired Techniques for Knowledge Discovery and Data Mining* (pp. 204-220).

www.irma-international.org/chapter/application-of-artificial-neural-network-and-genetic-programming-in-civil-engineering/110460