

Chapter 10

Implementation and Testing Details of Document Classification

ABSTRACT

It is trivial to achieve a recall of 100% by returning all documents in response to any query. Therefore, recall alone is not enough, but one needs to measure the number of non-relevant, for example by computing the precision. The analysis was performed for 30 documents to ensure the stability of precision and recall values. It is observed that the precision of large documents is less than a moderate length document, in the sense that some unimportant keywords get extracted. The reason for this may be attributed to the frequent occurrence and its unimportant role in the sentence.

SYSTEM TESTING

Reuters Data Set

Researchers have used benchmark data, such as the Reuters- 21578 corpus of newswire test collection (Sholom M. W., Indurkha, N., Zhang, T. and Damerau, F. 2010), to measure advances in automated text classification. We performed testing of our system using a sample of the same.

DOI: 10.4018/978-1-7998-3772-5.ch010

Modules of Execution

1. Document Entry
2. Stop Word removal
3. Stemming
4. Keyword generation
5. Document Classification

Document Entry

Doc_id : DOC1

Doc_content :

Table 1. Words after tokenization

hard
problem
text
classification
aspects
potential
solution
keyword
extraction
maximal
frequent
item
set
used
attributes
mining
association
rules
basis
measuring
similarity
new
documents
existing
association rules
issue
keyword
extraction
text
collection
emerging
research
filed
promotes
maximal
frequent
item
set
generation

9 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/implementation-and-testing-details-of-document-classification/268469

Related Content

Exploring Similarities Across High-Dimensional Datasets

Karlton Sequeira and Mohammed J. Zaki (2007). *Research and Trends in Data Mining Technologies and Applications* (pp. 53-84).

www.irma-international.org/chapter/exploring-similarities-across-high-dimensional/28421

Change Detection in Large Evolving Networks

Josephine M. Namayanja and Vandana P. Janeja (2019). *International Journal of Data Warehousing and Mining* (pp. 62-79).

www.irma-international.org/article/change-detection-in-large-evolving-networks/225807

Updating the Built Prelarge Fast Updated Sequential Pattern Trees with Sequence Modification

Jerry Chun-Wei Lin, Wensheng Gan, Tzung-Pei Hong and Jingliang Zhang (2015). *International Journal of Data Warehousing and Mining* (pp. 1-22).

www.irma-international.org/article/updating-the-built-prelarge-fast-updated-sequential-pattern-trees-with-sequence-modification/122513

The Healthcare Cost Dilemma: What Health Insurance Companies Can Do to Mitigate Unsustainable Premium Increases

Russ Danstrom and Jeff Nicola (2004). *Managing Data Mining: Advice from Experts* (pp. 124-145).

www.irma-international.org/chapter/healthcare-cost-dilemma/24782

A Fuzzy Portfolio Model With Cardinality Constraints Based on Differential Evolution Algorithms

JianDong He (2024). *International Journal of Data Warehousing and Mining* (pp. 1-14).

www.irma-international.org/article/a-fuzzy-portfolio-model-with-cardinality-constraints-based-on-differential-evolution-algorithms/341268