

Chapter 2

Grouping of Questions From a Question Bank Using Partition–Based Clustering

ABSTRACT

During automatic test paper generation, it is necessary to detect percentage of similarity among questions and thereby avoid repetition of questions. In order to detect repeated questions, the authors have designed and implemented a similarity matrix-based grouping algorithm. Grouping algorithms are widely used in multidisciplinary fields such as data mining, image analysis, and bioinformatics. This chapter proposes the use of grouping strategy-based partition algorithm for clustering the questions in a question bank. It includes a new approach for computing the question similarity matrix and use of the matrix in clustering the questions. The grouping algorithm extracts n module-wise questions, compute $n \times n$ similarity matrix by performing $n \times (n-1)/2$ pair-wise question vector comparisons, and uses the matrix in formulating question clusters. Grouping algorithm has been found efficient in reducing the best-case time complexity, $O(n \times (n-1)/2 \log n)$ of hierarchical approach to $O(n \times (n-1)/2)$.

TERMINOLOGY USED

The terminology used is presented in the table below -

Table 1. Terminology used for question clustering

Term	Meaning
Subject (S)	S is a subject/paper offered in different semesters of a course.
Modules/Units	For each subject, there is a university pre-scribed syllabus which consists of different modules/units.
Question Bank (QB)	QB is a database which stores module wise questions with its details such as question- no, question-content, question-type, question- marks and question-answer-time
Q	Q is the total number of questions stored under a module
t_i	t_i refers to the total number of questions in which term i appears
$freq_{ij}$	$freq_{ij}$ is the frequency of term i in question j
maximum frequency ($max\ freq_{ij}$)	$max\ freq_{ij}$ is the maximum frequency of a term in question j
term frequency (tf_{ij})	tf_{ij} refers to the importance of a term i in question j. It is calculated using the formula: $tf_{ij} = freq_{ij}/max\ freq_{ij}$
Inverse Document Frequency (idf_i)	idf_i refers to the discriminating power of term i and is calculated as: $idf_i = \log_2(Q/t_i)$
tf-idf weighting (W_{ij})	It is a weighting scheme to determine weight of a term in a question. It is calculated using the formula: $W_{ij} = tf_{ij} \times idf_i$
Question-Term-Set, T_i (question q_i)	A set of terms extracted from each question by performing its tokenization, stop word removal, taxonomy verb removal and stemming
Theshold, δ	User input threshold value to find the similarity

2. 2 PARTITION-BASED GROUPING ALGORITHMS FOR QUESTION CLUSTER FORMULATION

The similarity matrix computation has been carried out by using matrix representation of vectors which is a natural extension of existing Vector Space Model (VSM) (Jing, L., Ng, M. K.& Huang, J. Z. 2010; Turney, P. D., Pantel, P. 2010; Wong, S. K. M. and Raghavan, V.V. 1984). VSM is a popular information retrieval system implementation which facilitates representation of a set of documents as vectors in term space. Similarity matrix generates its

12 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/grouping-of-questions-from-a-question-bank-using-partition-based-clustering/268461

Related Content

A Clustering Rule Based Approach for Classification Problems

Philicity K. Williams, Caio V. Soares and Juan E. Gilbert (2012). *International Journal of Data Warehousing and Mining* (pp. 1-23).

www.irma-international.org/article/clustering-rule-based-approach-classification/61422

Measuring Semantic-Based Structural Similarity in Multi-Relational Networks

Yunchuan Sun, Rongfang Bie and Junsheng Zhang (2016). *International Journal of Data Warehousing and Mining* (pp. 20-33).

www.irma-international.org/article/measuring-semantic-based-structural-similarity-in-multi-relational-networks/143713

Sequential Patterns Postprocessing for Structural Relation Patterns Mining

Jing Lu, Weiru Chen, Osei Adjei and Malcolm Keech (2008). *International Journal of Data Warehousing and Mining* (pp. 71-89).

www.irma-international.org/article/sequential-patterns-postprocessing-structural-relation/1814

A Survey of Selected Software Technologies for Text Mining

Richard S. Segall (2009). *Handbook of Research on Text and Web Mining Technologies* (pp. 766-784).

www.irma-international.org/chapter/survey-selected-software-technologies-text/21757

The Application of an Integrated Behavioral Activity-Travel Simulation Model for Pricing Policy Analysis

Karthik C. Konduri, Ram M. Pendyala, Daehyun You, Yi-Chang Chiu, Mark Hickman, Hyunsoo Noh, Paul Waddell, Liming Wang and Brian Gardner (2014). *Data Science and Simulation in Transportation Research* (pp. 86-102).

www.irma-international.org/chapter/the-application-of-an-integrated-behavioral-activity-travel-simulation-model-for-pricing-policy-analysis/90067