

Chapter 1

Fuzzy-Based Querying Approach for Multidimensional Big Data Quality Assessment

Pradheep Kumar K.
BITS Pilani, India

Venkata Subramanian D.
Hindustan Institute of Technology & Science, India

ABSTRACT

This paper is intended to design a fuzzy based approach to assess standards and quality of big data. It also serves as a platform to organizations that intend to migrate their existing database environment to big data environment. Data is assessed using a multidimensional approach based on quality factors like accuracy, completeness, reliability, usability, etc. These factors are analysed by constructing decision trees to identify the quality aspects which need to be improved. In this work fuzzy queries have been designed. The queries are grouped as sets namely Excellent, Optimal, Fair and Hybrid. Based on the fuzzy data sets formed and the query compatibility index, a query set is chosen. A data set that has a very high degree of membership is assigned a fair query set. A data set with a medium degree of membership is assigned a optimal query set. A data set that has a lesser degree of membership is assigned a Excellent query set. A data set which needs a combination of queries of all the above is assigned a hybrid query set. The fuzzy query based approach reduces the query compatibility index by 36%, compared to a normal query set approach.

INTRODUCTION

In today's world with an increase in the amount of data processing and information requirement it is essential to develop strategies to effectively manage and assess the data for essential quality checks. The database forms the basis of day to day decisions taken by the organization. Data obtained from employees need to be periodically updated for effective utilization. In this work, an attempt has been made to assess data quality based on certain measures or parameters like Accuracy, Completeness, Reliability, Usability, etc as discussed by Pradheep et al in (2014). Based on these parameters the data set is queried to assess

DOI: 10.4018/978-1-5225-1008-6.ch001

the effectiveness of attributes like accuracy, usability, reliability, timeliness, etc. The parameters or quality factors such as Accuracy, completeness, etc are further subdivided into minor factors. Accuracy is subdivided into Syntactic and Semantic accuracy as explained by Pradheep et al (2014). The sub factors in turn have a parameter which is a measure and this parameter has an acceptable set of values.

To assess the effectiveness of the big data, a model is constructed for a Knowledge Management System which is a Multi-Dimensional Framework for quality checks. The different attributes are each modeled as a decision tree. The combination of all these form a Decision forest tree. A model for the decision forest tree was proposed by Criminisi et al in (2011). The decision forest tree model was a probabilistic model based on classification and regression analysis. The data under consideration could be textual, video, photographs to form a random forest of decision trees. To analyse this data effectively for information several data mining techniques were proposed. Berendt and Preibusch (2014) have proposed several techniques to extract data from databases based on the choice of attributes. Another technique which is map and reduce technique for large data sets had been proposed by Doukeridis and Norvag in (2014). A large number of data visualizing techniques have been explained in this regard by Doukeridis and Norvag (2014), Venkat et al (2011), Gorodov et al (2013), Serban et al (2013), Shamsi et al (2013), Jennex and Olman (2003), Evans et al (2013) and Banerjee et al (2014).

A data dictionary needs to be available which acts as a repository for storing the data. The data dictionary would contain metadata of the data related to the nature, type, volume, etc. Data integrity is another feature which decides on the reliability of data. Data access should be provided according to the role based privileges. This is done based on access privileges and functional aspects. The access privileges may vary from time to time based on the effectiveness of the queries which are also assessed to ensure minimal processing time and memory. Based on this approach the queries are classified into sets. Based on the decision tree analysis carried out the entire data is partitioned into smaller datasets. The size of the dataset may vary arbitrarily based on data volume and processing speed of the database.

The datasets are discretised based on the degree of membership. The queries are classified based on the query processing time and memory. A data set that has a very high degree of membership would have data points very closely spaced. This type of data set has a scattering distance which is negligible. The distance is measured by Euclidean method. In other words the data points would be clustered and would belong to only a single fuzzy set. Hence it would be ideal to use a query set which has a low query compatible index. This would not have a very high impact on the query cost, thereby optimizing the same. When the degree of membership is partial, we would need a slightly more efficient query set to work on this type of data set. Here the data set is spaced a bit and may belong to two different fuzzy sets. To accomplish this, the query would need to work on both fuzzy sets to return precise information. In such cases the data points are separated by a finite scattering distance. Here there is no clustering. For these data sets the query set should have an average or optimal query compatibility index.

When the degree of membership of the data is very low, the data points are widely scattered with a very large scattering distance. Here there is no question of clustering and the data points are distributed among multiple fuzzy sets. This would imply a very large scattering distance. To extract precise information, it is essential to assign highly efficient query set with a very high query compatibility index. A Mamdani's fuzzy inference engine is used as the inputs are non-linear and are provided as fuzzy sets. In this work a triangular membership function has been used to determine the degree of membership. The entire approach aims at ensuring no compromise in performance and optimal assignment of datasets with query sets.

21 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/fuzzy-based-querying-approach-for-multidimensional-big-data-quality-assessment/169480

Related Content

Big Data Virtualization and Visualization: On the Cloud

Muhammad Adeel (2017). *Decision Management: Concepts, Methodologies, Tools, and Applications* (pp. 1436-1452).

www.irma-international.org/chapter/big-data-virtualization-and-visualization/176813

Performance Assessment of R&D-Intensive Manufacturing Companies on Dynamic Capabilities

Mohammadyasser Darvizeh, Jian-Bo Yang and Stephen Eldridge (2020). *International Journal of Strategic Decision Sciences* (pp. 1-23).

www.irma-international.org/article/performance-assessment-of-rd-intensive-manufacturing-companies-on-dynamic-capabilities/269686

A Decision Support System for Managing Demand-Driven Collection Development in University Digital Libraries

Mohamed Hemili, Mohamed Ridda Laouar and Sean B. Eom (2021). *Research Anthology on Decision Support Systems and Decision Management in Healthcare, Business, and Engineering* (pp. 901-919).

www.irma-international.org/chapter/a-decision-support-system-for-managing-demand-driven-collection-development-in-university-digital-libraries/282622

Towards a Framework for the Measurement and Reduction of User-Perceivable Complexity of Group Decision-Making Methods

Andrej Bregar (2014). *International Journal of Decision Support System Technology* (pp. 21-45).

www.irma-international.org/article/towards-a-framework-for-the-measurement-and-reduction-of-user-perceivable-complexity-of-group-decision-making-methods/123993

The Development and Measurement of a Customer Satisfaction Index (E-CSI) in Electronic Banking: An Application to the Central Vietnam Region

Le Van Huy, Pham Long, Aidan O'Connor and Pham Dinh Tuyen (2017). *International Journal of Strategic Decision Sciences* (pp. 45-58).

www.irma-international.org/article/the-development-and-measurement-of-a-customer-satisfaction-index-e-csi-in-electronic-banking/189077