

Automatic Field Identification for Tabular Data without Headers

Venkat Giridhar Mareedu¹ Keerthi Kant Raghunathan²

¹MVSR Engineering College ²Osmania University India

Abstract— finding useful patterns from large sets of data has brought new possibilities and insights which are unnoticeable to the human eye. But this is all possible only when we know what data it is i.e., when we have information about the metadata of the database under consideration. This paper aims to automatically identify table headers and their datatypes without any human intervention. Upon identification, this paper also implements a strategy to automatically import the data into a big data tool, SAS, to conduct further data analytics. These strategies are perfectly consistent even with a table with headers as it saves the time of the user who would otherwise manually map the fields.

Key words: Big Data, Data Mining, SAS, Artificial Intelligence, Clustering

I. INTRODUCTION

With the fast evolving world, the amount of data that is generated every day is increasing exponentially. This has generated an urge to process data in real-time. Real-time data processing enables for faster analysis and responses which could have sizeable impact on the operation in hand. Methods and various strategies have been proposed to automate the task of importing the data and performing the requisite analysis. Though largely successful, there are issues which have prevented the system from being a full automated one.

Any hindrance or presence of abnormalities in handling large chunks of data can create potential bottle necks. One such problem largely persistent is that of the unavailability of metadata information of a table. Metadata of a data is the information about the data. The no. of fields, the information that those fields contain, the types of values the fields contain are some of the important information that a metadata of a table contains. Generally when such data is not available, human intervention is required to map the fields and import the data so that the analysis can proceed further. With the use of several data mining techniques, this paper aims to automate this step of manually mapping the fields and importing the data.

Initially, the data to be imported is sampled in a random fashion and a subset is created. This subset is then processed each column at a time by finding the nearest column to it from a pre-defined dataset which has properly defined clusters wherein a cluster corresponds to a respective column. Ideally, by the end of the process, each column to be mapped would be a part of a single cluster. If not, the process is then repeated by picking another larger subset from the data to be imported. This paper mainly focuses on identification of columns from a student database. Note that the process discussed is application specific and can be adapted accordingly as the dataset to be imported has to match at least a minimum number of columns in the pre-defined clustered dataset.

II. METHODOLOGY

A. Data Mining

Let us refer the dataset for which the headers have to be identified as Dataset A and the dataset which is used train so as to form pre-defined clusters as Dataset B.

Initially, before the dataset A is input, we create clusters of clearly defined boundaries using Dataset B for which we have the metadata and are able to form clusters accordingly. Crucial information such as the type of the fields and the maximum length of the fields is captured so as to help in creating the schema for datasets which are to be input later. There are many techniques which can be used to be form the requisite clusters. But the results drawn from the following two algorithms have been significantly better than the rest for the student dataset used:

- k-Means Algorithm
- DBScan Algorithm

B. k-Means Algorithm

It is an iterative algorithm in which objects (a column value of each record) are moved among a set of clusters until the desired set is achieved. The algorithm is built on the concept of user specified input parameter (k) which is the number of clusters to be formed. This outperforms DBScan with its high precision due to fact that the number of clusters are specified beforehand which is not the case in DBScan algorithm. As we have the metadata of Dataset B, this algorithm is more suitable for our scenario. [1]

C. DBScan Algorithm

Density-based clustering locates regions of high density that are separated from one another by regions of low density i.e., it creates the clusters with a minimum size and density. This contrary to k-means, handles the outlier problem by ensuring that an outlier will not create a cluster. But as the number of columns/clusters are not specified, the algorithm fails to classify clusters with fine margins. [1]

D. Classification Techniques

Once the clusters are formed using the dataset B, the dataset A can be input to the model to determine headers and other Meta information regarding the data. There are several classifying techniques which can be applied to determine which column corresponds to which but the combination of multiple classifying techniques gives more assured results but at the cost of run time. The following combination of classification techniques works best for the student database:

1) K-Nearest Neighbors Classifier

Finding out the k nearest neighbors using simple Euclidean distance formula. The mean of all the clustered tuples is calculated and is compared against the individual tuple values to determine which column falls under which cluster. For character fields, attribute values can be used to determine the closeness among the fields under comparison. [3]

REFERENCES

- [1] Data Mining: Concepts and Techniques (The Morgan Kaufmann Series in Data Management Systems) by Jiawei Han, Micheline Kamber, Jian Pei Professor
- [2] <http://poi.apache.org/spreadsheet/quick-guide.html#CreateCells>
- [3] https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm
- [4] Y. LeCun, L. Bottou, G. Orr, and K. Muller, "Efficient BackProp," in Neural Networks: Tricks of the trade, (G. Orr and Muller K., eds.), 1998.

