

FMCleaner: Automatic detect and repair data error using Functional Dependencies and Machine Learning algorithm

Jesmeen M. Z. H.^{1,*}, J. Hossen¹, K. Tawsif¹, Md. Armanur Rahman¹, Shohel Sayeed²

¹) Faculty of Engineering and Technology, Multimedia University, Melaka, 75450, Malaysia

²) Faculty of Information Science & Technology, Multimedia University, Melaka, 75450, Malaysia

*Corresponding e-mail: jesmeen.online@gmail.com

Keywords: Data cleaning; Machine Learning; Integrity constraints

ABSTRACT – Data analytics (DA) technology beneficial for analyzing and predicting, and it also helps to make decisions. In the process of DA one of the challenges are detection and repairing dirty data, where failure to do so can result in incorrect analytics and defective decisions. In this paper, an FMCleaner is proposed and implemented to automatically detect (such as Integrity Constraints) and repair data error. The developed system evaluated by using a test dataset with an accuracy of around 95%.

1. INTRODUCTION

Data quality became a major concern for many organizations. As stated by [1] “Data cleaning, as an essential aspect of quality assurance and a determinant of study validity, should not be an exception”. Few reasons for producing an error is dataset are missing data, different formats (such as date format), replicated entered data, typos, outlier data and violating business rules.

For the past few years, on emerging new data, industries and academics fields become worried about data cleaning [2–4]. For different format of data required different scripting. However, Machine Learning (ML) have ability to develop an automatic system for cleaning data [5]. Whereas, integrity constraints (ICs) helps to elaborate rules of data cleaning

Contributions. The main contributions in this work are to present current technique for progressive detecting data error and cleaning data automatically using ML as highlighted in figure 1. Moreover, evaluation of the system was executed using four data sets obtained from UCI repository and Malaysia TM Company and it results in better prediction accuracy. Finally, in this paper, result and discussion one of the dataset outcomes discussed.

2. DATA CLEANING TECHNIQUES

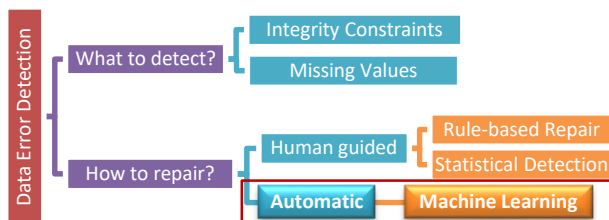


Figure 1 Proposed system data error detection

As shown in figure 1, the system proposed contains two main phases: detection and repairing.

2.1 Error Detection

Detecting IC depending on first-order logic

including Functional Dependencies (FDs) to understand the rules of captured data quality. While missing values can also produce issues in analytical process. Designing ICs can obtain manually, but we had proposed a system, which will execute FDs automatically.

2.2 Error Repair

The proposed approaches for repairing dirty data involving ML for predicting missing values, rather than setting missing values in statistic method (such as mean/median). Other methods of cleaning data involved human guidance to confirm the fixes, suggestion to fix or to select the best ML models to process automatic repairing decisions [6].

3. PROBLEM FORMALIZATION

Set of Rules (R) obtained to use for identifying a set of inconsistent values (such as errors due to typos) for each categorical fields.

Pseudocode to discover FD using Pruning Algorithm

Inputs: D is the training set

Col is columns of the dataset

Output: Function Dependencies for input dataset

rulesFD={}

C0:={}//Empty list of columns

C1:= Col//is a copy of listofcolumns

c:=1

C:=[C0,C1]

while Cc!=0

 compute_dependencies(Cc, C1)

 PRUNE(Cc)//helps to Reduce column combinations

 temp:= generate_next_level(Cc)

 Cc+1.append(temp)

 c:=c+1

end

Trained ML model used for numerical and non-numerical fields to predict missing values. In this case, we assumed the training set is already clean (D_{clean}) manually by inspecting the data as far as possible. The system contains list of function $F = \{f_1, \dots, f_n\}$ to repair detected data issues. Considering, D_{Dirty} to be set of containing data error and repair with selected function from F.

The implemented algorithm is as follows:

1. Let C_{null} be the missing value columns, Detect C_{null} containing missing values using F_{NAN} (Function to detect is column contains null/NaN/?)
2. Obtaining, D_{clean} (the dataset cleaned previously), and let C_{clean} is list of labels/features/attributes labels for each record containing clean data

